

UNCLASSIFIED

AD NUMBER
AD024361
NEW LIMITATION CHANGE
TO Approved for public release, distribution unlimited
FROM Distribution authorized to DoD only; Administrative/Operational Use; OCT 1952. Other requests shall be referred to Bureau of Aeronautics, Department of the Navy, Washington, DC 20350. Pre-dates formal DoD distribution statements. Treat as DoD only.
AUTHORITY
NAVAIR ltr dtd 6 Jun 1978

THIS PAGE IS UNCLASSIFIED

AD No. 24 361

ASTIA FILE COPY

BuAer Report AE-61-4

Fundamentals of Design
of Piloted Aircraft
Flight Control Systems

Volume I

**METHODS OF ANALYSIS AND
SYNTHESIS OF PILOTED AIRCRAFT
FLIGHT CONTROL SYSTEMS**

PUBLISHED BY DIRECTION OF
THE CHIEF OF THE BUREAU OF AERONAUTICS

**Best
Available
Copy**

BuAer Report AE-61-4

**Fundamentals of Design
of Piloted Aircraft
Flight Control Systems**

Volume I

**METHODS OF ANALYSIS AND
SYNTHESIS OF PILOTED AIRCRAFT
FLIGHT CONTROL SYSTEMS**

**PUBLISHED BY DIRECTION OF
THE CHIEF OF THE BUREAU OF AERONAUTICS**

OCTOBER 1952

PREFACE

This volume has been written under BuAer Contract NOs 51-514(c) to provide engineers with analytical and other techniques basic to the unified approach to problems of aircraft control system design.

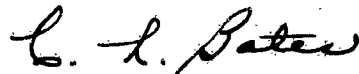
A large portion of the volume is a codification of existing techniques and material appearing in textbooks and published papers. However, a certain amount of new material is also presented for the first time in published form. Bibliographies covering the major source material are appended to each chapter.

This volume has been written from the point of view that the basic approach to control systems problems is of necessity through the transfer function. The various ways of dealing with such problems in practice are essentially means of getting various degrees of approximations to the transient solution of the equations of motion from which the transfer functions are derived. Since the object of prime interest in control and servomechanisms work is the transient behavior of the system under consideration, it is felt that this approach will provide the control systems engineer with a relatively new codifying concept with which to attack his problems.

The authors are indebted to many individuals and companies who have aided or influenced this volume either directly or indirectly, and particularly to the Bureau of Aeronautics of the United States Navy. Special appreciation is due to Mr. L. M. Chatlier, Mr. R. A. Benneche, and Mr. J. Folse, of BuAer, whose foresight and continued interest have made this project possible. Special mention should be given to Juanita Zimmerman, Betty Harsey, Elias Moness, F. B. Bacus, and James Jones of the Northrop Servomechanisms Section for their untiring efforts in preparing the manuscript for publication; also to K. B. Tuttle, who was charged with the responsibility of coordinating all the individual efforts involved. The arduous task of writing this book has been made a more pleasant one by the continuing interest and able assistance of all those mentioned above and the entire Northrop Engineering Division.



D. T. McRuer, Supervisor
Servomechanisms Section



C. L. Bates, Director
Mechanical Design Department

EDITORIAL BOARD

R. J. Kulda	K. B. Tuttle
E. Moness	R. Zacharias
D. T. McRuer	

CONTRIBUTING AUTHORS

B. C. Axley	D. T. McRuer
R. G. Halliday	E. Moness
O. Imai	K. B. Tuttle
R. J. Kulda	R. Zacharias
L. H. Lyons	

IMPORTANT NOTE

This volume was written by and for engineers and scientists who are concerned with the analysis and synthesis of piloted aircraft flight control systems. The Bureau of Aeronautics undertook the sponsorship of this project when it became apparent that many significant advances were being made in this extremely technical field and that the presentation and dissemination of information concerning such advances would be of benefit to the Services, to the airframe companies, and to the individuals concerned.

A contract for collecting, codifying, and presenting this scattered material was awarded to Northrop Aircraft, Inc., and the present basic volume represents the results of these efforts.

The need for such a volume as this is obvious to those working in the field. It is equally apparent that the rapid changes and refinements in the techniques used make it essential that new material be added as it becomes available. The best way of maintaining and improving the usefulness of this volume is therefore by frequent revisions to keep it as complete and as up-to-date as possible.

For these reasons, the Bureau of Aeronautics solicits suggestions for revisions and additions from those who make use of the volume. In some cases, these suggestions might be simply that the wording of a paragraph be changed for clarification; in other cases, whole sections outlining new techniques might be submitted.

Each suggestion will be acknowledged and will receive careful study. For those which are approved, revision pages will be prepared and distributed. Each of these will contain notations as necessary to give full credit to the person and organization responsible.

This cooperation on the part of the readers of this volume is vital. Suggestions forwarded to the Chief, Bureau of Aeronautics, (Attention AE-612), Washington 25, D. C., will be most welcome.

L. M. Chatter
Head, Actuating & Flight Controls Systems Section
Airborne Equipment Division
Bureau of Aeronautics

TABLE OF CONTENTS

CHAPTER I:	GENERAL CONSIDERATIONS	I-1
CHAPTER II:	FUNDAMENTAL CONCEPTS	
Section 1-	Introduction	II-1
Section 2-	Mathematical Models	II-1
(a)	General	II-1
(b)	Linearization	II-2
Section 3-	Block Diagrams and Transfer Function	II-5
(a)	The Block Diagram	II-5
(b)	Block Diagram Algebra	II-8
(c)	Equivalent Block Diagram	II-16
(d)	The Transfer Function	II-18
(e)	Graphical Forms	II-26
Section 4-	Servomechanisms	II-37
CHAPTER III:	ANALYSIS	
Section 1-	Introduction	III-1
Section 2-	The Generalized Nyquist Method	III-2
(a)	Closed-Loop s -Plane	III-2
(b)	The Mapping Theorem	III-2
(c)	The Conventional Nyquist Criterion	III-3
(d)	Specified Minimum Damping and Damping Ratio	III-5
Section 3-	The Open Loop-Closed Loop Logarithmic Method	III-10
(a)	General	III-10
(b)	Relations Between Open and Closed-Loop Transfer Functions-Charts	III-12
(c)	Relationship Between Open and Closed-Loop Transfer Functions--by Approximation	III-14
(d)	Development of Analytical Form of Closed-Loop Transfer Function from Graphical Representation	III-16

Section 4-	The Root Locus Method	III-20
(a)	Introduction	III-20
(b)	Basic Principles	III-20
(c)	Constructing the Locus	III-23
(d)	The Root Locus Plotter	III-27
(e)	Applications	III-32
Section 5-	Special Cases — Gain Margin, Phase Margin, Maximum Magnification Ratio	III-33
CHAPTER IV: SYNTHESIS		
Section 1-	Introduction	IV-1
(a)	General	IV-1
(b)	System Design Procedure	IV-1
(c)	Equalizer Synthesis	IV-2
Section 2-	Specification of Servo System Performance	IV-2
Section 3-	Equalization	IV-4
CHAPTER V: OPTIMUM SYNTHESIS METHODS		
Section 1-	Introduction	V-1
Section 2-	Fixed System Optimization	
Section 3-	Variable System Optimization	
CHAPTER VI: NON-LINEARITIES		
Section 1-	Introduction	VI-1
Section 2-	Continuous Non-Linearities	VI-2
Section 3-	Discontinuous Non-Linearities	VI-3
(a)	General	VI-3
(b)	Small Discontinuities	VI-4
(c)	Phase Plane	VI-12
CHAPTER VII: MACHINE METHODS		
Section 1-	Introduction	VII-1
Section 2-	Need for Machine Methods	VII-1
Section 3-	Available Methods of Automatic Computation	VII-2
(a)	Digital Computers	VII-2

(b)	Analog Computers	VII-3
Section 4-	Evaluation of Machine Methods for Control System Work	VII-4
(a)	Set-up Time	VII-4
(b)	Variation of Parameters	VII-4
(c)	Form of Data	VII-5
(d)	Simulation	VII-5
(e)	Accuracy	VII-5
 CHAPTER VIII: THE ANALOG COMPUTER		
Section 1-	Introduction	VIII-1
Section 2-	"Ideal" D.C. Amplifier Operation	VIII-1
Section 3-	Non-Ideal Operational Amplifiers	VIII-3
Section 4-	Effects of Non-Ideal Operation	VIII-5
Section 5-	Solution of Differential Equations	VIII-6
Section 6-	Simulation of Non-Linearities	VIII-10
(a)	Introduction	VIII-10
(b)	Coulomb Friction	VIII-11
(c)	Spring Preload	VIII-12
(d)	Threshold	VIII-12
(e)	Limiting	VIII-12
(f)	Hysteresis	VIII-12
(g)	Continuous Non-Linearities	VIII-13
 APPENDIX: MATHEMATICAL BACKGROUND		
Section A	Tables, Charts and Graphs	A-1
Section AI	Roots of Algebraic Equations	A-55
Section AII	The Routh-Hurwitz Stability Criterion	A-59
Section AIII	Function of a Complex Variable	A-60
Section AIV	Mapping	A-66
(a)	Introduction	A-66
(b)	Methods of Designating Points, Lines and Areas in the Z-Plane	A-66
(c)	Mapping	A-68
Section AV	Factoring Polynomials by Servo Analysis Methods	A-72
GLOSSARY		B-1

CHAPTER I

GENERAL CONSIDERATIONS

The purpose of this volume is to present the systems engineer with essential mathematical tools for solving problems arising in piloted aircraft control system design. This introductory chapter is devoted to surveying some of the general principles involved and to define certain important terms.

As far as his degree of control over the airplane is concerned, a pilot may be said to engage in two types of flight activity in the course of a mission. The first of these, usually called navigation, does not require great precision of control but enables the pilot to fly the airplane to some point at which he begins the second type of flight activity characterized by the much greater control needed. During this second type of flight, the pilot must direct the airplane in some precise geometrical relationship to an object on the ground or in the air. The object on the ground may perhaps be the landing field where the pilot will terminate the flight or a target which is to be bombed; the object in the air may perhaps be an enemy aircraft which is to be destroyed or a friendly aircraft with which the pilot is to fly in completing a mission. Depending upon the particular situation, it can be seen that this precision type of flight is required in such activities as landing, tracking, gunlaying, and bombing runs. The essential

difference between the two types is one of precision of control maintained by the pilot; the essential likeness is that the pilot is primarily interested in directing or commanding the airplane to assume a certain orientation in space.

The pilot uses the cockpit control devices for two functions: to command and to stabilize the airplane. The difference between these two functions can be clarified by a simple but common illustration:

A pilot takes an airplane aloft and at a certain flight condition attempts to establish a certain rate of climb. He pulls back on the stick and holds it steady. In response to this command, the airplane smoothly takes up a steady climb. However, the climb rate is not precisely what the pilot wants, so he applies a slight forward force to the stick. The airplane then pitches forward, and the rate of climb is reduced. But it is reduced too much, so the pilot applies back pressure again, and again the airplane climbs too fast. This procedure continues: The pilot adjusts the stick forward and backward continuously, and the airplane pitches up and down without ever settling down to the desired rate of climb.

In this example, the pilot was unable to command a rate of climb and get exactly what he wanted. He had to jockey the stick back and forth, climbing first too fast, then too slowly. The average rate of climb was probably satisfactory, but the pilot was compelled to work very hard to maintain it. If he has been able to stabilize at the desired rate of climb after perhaps one or two small corrections with the stick, he would have been satisfied; but he was not satisfied with the oscillating response he actually achieved.

However, performance was unsatisfactory only if the pilot insisted on getting exactly the rate of climb he had initially decided upon. As long as he placed the stick in a certain position and held it there (command input), the airplane flew smoothly with no tendency to oscillate; its flight behavior was stable. It was only when he attempted to make certain fine adjustments that the oscillation occurred. The airplane by itself was stable. What was unstable was actually the pilot-airplane combination.

To solve a problem of this sort, one must study three factors: a) the pilot, to determine how he responds to certain stimuli, such as pitch angle, normal acceleration, and stick force; b) the airplane, to determine how it responds to certain control movements and air forces, such as backward and forward stick deflection, lift changes, and pitching moment changes; and c) the pilot-airplane combination as a system, to determine how they will interact with each other.

When this has been accomplished, changes can be made in certain elements to produce a better system.

This volume presents techniques suitable for handling feedback control systems problems of types of practical importance to engineers concerned with the design of flight control systems, of which the one considered above is an example. These are the methods of analysis: methods of determining how a system or elements of a system behave when subjected to command inputs or to external disturbances. The volume goes beyond analysis, however, and consider methods of synthesis, that is, procedures for determining the best way of selecting many elements and combining them into a system. Analysis takes a given system and determines its behavior, whereas synthesis creates a system which will behave according to a certain desired pattern.

In building up a control system, it is usually convenient at first, and often necessary, to "live with" a certain number of components. In controlling an airplane, for example, the airframe must usually be accepted without change because certain design parameters affecting its performance were determined by considerations other than control, such as landing speed and maximum gross weight. Those elements which are accepted without change are referred to as unalterable elements. All the other elements in the system which can be chosen or designed at will to obtain the desired performance are called alterable elements. In the example used above, certainly the pilot and most probably the airframe would be considered unalterable elements. The surface controls would be the alterable elements. Thus a problem in analysis can be considered

as one in which the behavior of the system is studied to indicate which of the alterable elements might be changed to produce satisfactory performance. On the other hand, synthesis is concerned with designing or choosing alterable elements which will enable the system to produce a desired performance.

The term system, as used in this volume, may be defined as any group of interacting entities required to account completely for the physics of a particular observed phenomenon. This definition includes systems as simple as a pendulum and as complex as an aircraft fire control system. For the pendulum, the interacting entities are gravity, the atmosphere, bearing friction, and the mass of the pendulum, as illustrated in figure 1-2. In an aircraft fire control system, there are too many elements to list here. However, a pertinent point is that many of the elements are themselves complex "sub-systems," such as a tracking radar system or a pilot-controls-airframe combination similar to the one described earlier in this chapter.

The purpose of this survey has been to present some general considerations relating to the design of piloted aircraft control systems: in subsequent chapters, the mathematical means for dealing with design problems are discussed.

This volume is directed to college graduate engineers. Consequently, the entire effort is toward a logical presentation of the methods of analysis and synthesis. This precludes digressions for the presentation of general background material. However, since it is expected that the engineers who use this book will have backgrounds in various branches of engineering, not all of which make use of the mathematics pertinent to controls analysis, an appendix is included which touches briefly on important topics. It is highly recommended that the reader scan the appendix so that he may acquaint himself with any unfamiliar material found there.

In addition to certain mathematical material, the appendix contains a rather extensive glossary of terms. In particular, many new definitions are given to old and familiar words. Although this may at first seem arbitrary, experience will show that once these definitions are well fixed in mind, the text can be followed with much less confusion than would exist had the definitions been left to chance. The reasons for this are simply that automatic control theory is relatively new and that there are a number of systems of nomenclature in use. If no standardization were settled upon for this book, no two readers would gather the same impressions from the text. The definitions adopted by the AIEE-ASME joint committee at the time this volume is published are used whenever possible. A list of these symbols and definitions is given in the glossary, together with those of whatever standards exist in fields where the AIEE-ASME presentation is not applicable.

The main body of text can be divided into three major divisions. The first division is simply chapter II in which fundamental background material relevant to the description of system behavior is presented. Chapters III, IV, and V form the next major division. These

chapters are concerned solely with the analysis and synthesis of linear systems. The third division consists of chapters VI, VII, and VIII and considers not only linear devices, but also non-linear systems. Chapters

VII and VIII are devoted to machine methods of handling analysis and synthesis problems. These chapters complete the volume.

CHAPTER II

FUNDAMENTAL CONCEPTS

SECTION I — INTRODUCTION

This chapter discusses the means of describing control systems in order to facilitate analysis and synthesis. The mathematical form used to describe a system and each of its components is called "the transfer function." The transfer function of a linear system completely specifies its dynamic performance in terms of the Laplace transform variable s and certain fundamental parameters such as natural frequency, time delay and damping ratio. Since control systems analysis and synthesis is carried on almost entirely in terms of transfer functions, it is important that its forms, meanings and various graphical representations be

firmly established.

The transfer function is essentially a "mathematical model" of a unit or system, and it is manipulated like a laboratory model in order to produce satisfactory performance.

Manipulation of the transfer function is greatly facilitated by a number of graphical aids that show relationships between the transfer functions of the individual components of a system, and promote visualization of the transfer functions themselves.

SECTION 2 — MATHEMATICAL MODELS

(a) GENERAL

Given the problem of designing a gear unit, an engineer proceeds by first considering the physical requirements, such as mechanical advantage, pitch diameter, and stresses. Bearing in mind such further factors as producibility and availability of material, he formulates his design. These ideas are then converted into working drawings, from which the gear is produced. Finally, the finished product is tested.

Essentially, the design procedure may be summarized as follows:

1. The requirements of the design are determined.
2. Calculations are performed in order to determine controlling parameters.
3. The best methods of meeting requirements are determined.
4. Working drawings are prepared.
5. The unit is constructed.
6. The completed design is tested.

A similar procedure is followed in designing control systems. However, the complexity of such systems makes a more extensive design procedure necessary. Following the formulation of the design, a rigorous analysis must be applied to answer these questions:

1. Is the basic concept sound?
2. How does the system perform?
3. How well does the system perform?
4. How can the system be improved?

These analyses can sometimes be done through laboratory experiments. As an example, consider the analysis of an electrical distribution network. The system is simulated in the laboratory by assuming that

lumped parameters approximate the distributed capacitances, inductances, and resistances. These parameters, along with variables such as loading conditions to represent disturbances, are varied and measurements of the responses are made. Because of the assumption of lumped parameters, the results are not a perfect representation of the physical system. However, they are the best obtainable under laboratory conditions and may yield a great deal of useful information.

The complex nature of most control systems often precludes the laboratory experiment method of analysis. However, one can come close to achieving the advantages of the laboratory experiment by performing the experiment on paper. Mathematical equations may be used to represent a physical system in much the same way as the laboratory model simulates it. When properly derived both the mathematical and laboratory models should reflect, within certain limits, the characteristics of the physical system being analyzed. What these limits are is determined by the assumptions used. Obviously, a mathematical model derived from unfounded assumptions will lead to grossly misleading conclusions. Consequently, assumptions must be made with great care. Precautions must be taken throughout the analysis not to lose sight of these assumptions and to realize that they limit the validity of the final conclusions.

In order to provide a reminder, it is wise to make a complete and well defined list of assumptions which control the analysis. This procedure yields additional benefits. For one thing, in order to express precisely what he has in mind, the engineer must think critically about his statements. For example, it is easy to make the mental assumption "no friction," but it requires more careful consideration to write down precisely

Chapter II Section 2

where there is no friction. In considering "where," the engineer often discovers that the assumption is unreasonable and hence changes his attack. Another additional benefit derived from the list of assumptions is realized when discrepancies occur between analytical and experimental data. Reference to the list usually points to the cause.

For these and many other reasons, it is highly recommended that all assumptions be listed before proceeding on any analysis problem.

(b) LINEARIZATION

One of the most important assumptions usually made in analysis concerns the "linearity" of the elements of the system. The system shown in figure II-1 is considered in some detail in section II-3 where the following equation is derived.

$$(II-1) \quad M \frac{d^2x}{dt^2} + B \frac{dx}{dt} + kx = Q(t)$$

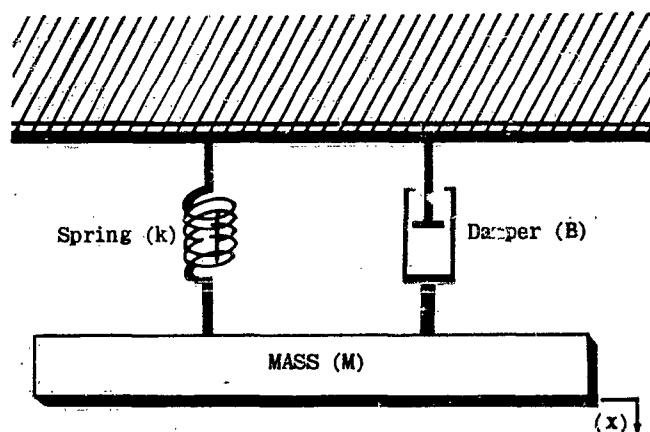


Figure II-1. Second Order System

Since x is the displacement of the mass, this equation says that the sum of the inertia force (Md^2x/dt^2), the damping force (Bdx/dt), and the spring force (kx), equals the driving force, $Q(t)$. In order to solve this equation conveniently, the following three assumptions are made:

1. The mass M is constant.
2. The damping factor B is constant.
3. The spring characteristic k is constant.

If these statements are true, (II-1) is of the form:

$$(II-2) \quad a_0 \frac{d^2x}{dt^2} + a_1 \frac{dx}{dt} + a_2 x = Q(t)$$

where a_0, a_1, \dots, a_2 are constants and $Q(t)$ is some function of time (t). Equations of the form of (II-2) are called linear differential equations with constant coefficients and are readily solved. However, none of the three assumptions is precisely correct. Consider, for example, the spring characteristic k . One form which this quantity may have is illustrated in figure II-2.

The essential point here is that the spring characteristic k is a function of the dependent variable x . If M and

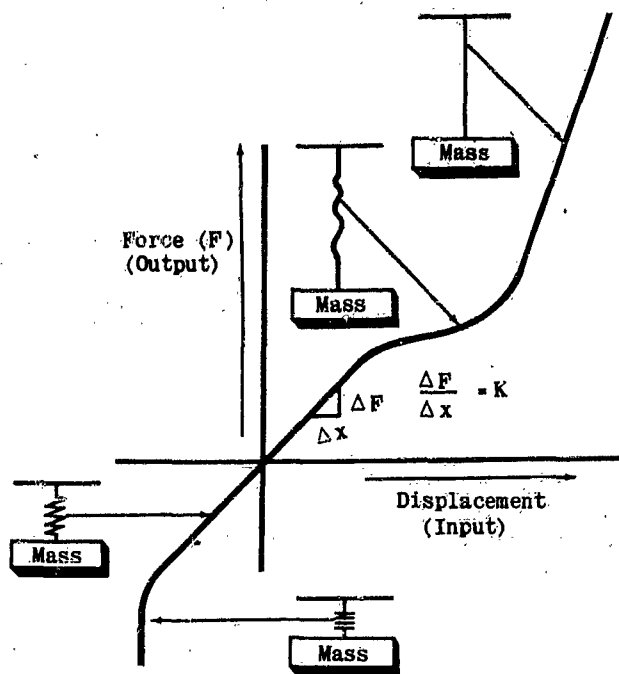


Figure II-2. Non-Linear Spring Constant

B vary in a similar fashion, (II-1) takes the form:

$$(II-3) \quad f_0(x) \frac{d^2x}{dt^2} + f_1(x) \frac{dx}{dt} + f_2(x) = Q(t)$$

This is known as a non-linear differential equation. In the simple example chosen it was assumed that the $f_i(x)$ were functions of x only. In general, however,

$$f_i(x) = f_i \left[x, \frac{dx}{dt}, \frac{d^2x}{dt^2}, \dots, \frac{d^nx}{dt^n} \right]$$

In most cases equations of the type of (II-3) are solved only by lengthy computations. However, in some cases, such as that represented by figure II-2, an additional assumption can be made that reduces the equation to the form of (II-2). Specifically, it can be assumed that the range of values of x is restricted so that the slope of the curve is constant. This is called the linear range of x .

It can be inferred from the discussion that one necessary characteristic of a linear system is that the static relationship between the input and the output is a straight line. Systems which do not satisfy this criterion are said to have non-linear static characteristics.

Non-linear static characteristics are divided into two main classifications:

1. Continuous non-linearities.
2. Discontinuous non-linearities.

Figure II-2 represents a continuous non-linearity. Figure II-3a represents a different type of non-linear static characteristic in that there is no straight line portion. The device represented might be a pressure transducer, figure II-3b, in which case the input is pressure and the output is volts. The curved line on the graph of input vs. output represents the actual static characteristic of the system. The dotted straight

lines represent several different possible linear approximations to the true curve. If the operating range of interest were large and symmetrical about the origin, the line xy would be used. If the operating range were small, the line uv would be used. If the range of operation centered about the point a , the straight line passing through that operating point would be used to represent the system static characteristic. Whenever the true characteristics of a system are approximated by a straight line in this manner, it is said that the system has been "linearized."

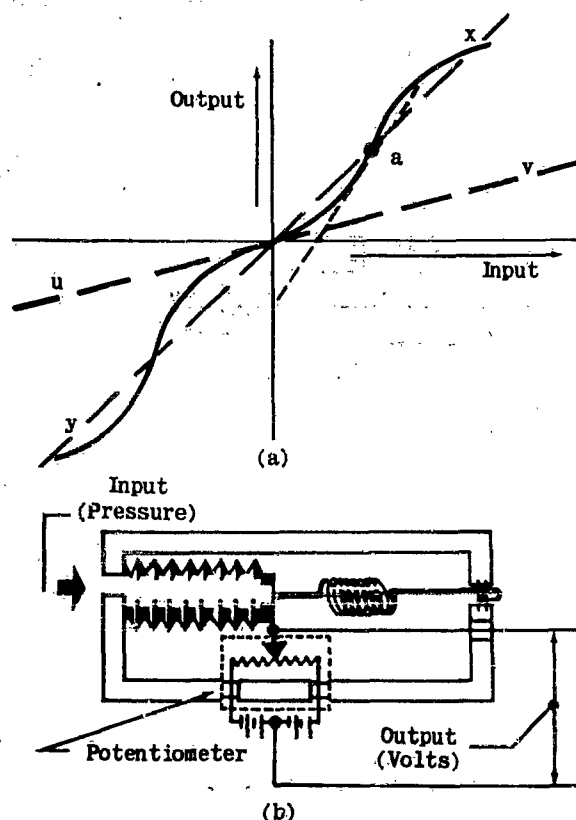


Figure II-3. Continuous Non-Linear System

Suppose that the pressure transducer were redesigned so that it had a static characteristic represented by the line xy . Next, assume that the potentiometer wiper attachment to the bellows became loose. The unit would then appear schematically as in figure II-4b. The bellows can now move through a certain distance without moving the potentiometer wiper. The result is a discontinuous static characteristic as shown in figure II-4a. This type of discontinuous non-linearity has as its essential feature a hysteresis loop and is called backlash. Note that it can be referred to in terms of the output or input.

The only way to linearize a system of this type is to assume that the backlash is negligible. This is equivalent to saying that the range of operation of greatest concern is large compared to the backlash. This situation is illustrated in figure II-5.

Figure II-6 illustrates some other important discontinuous non-linearities.

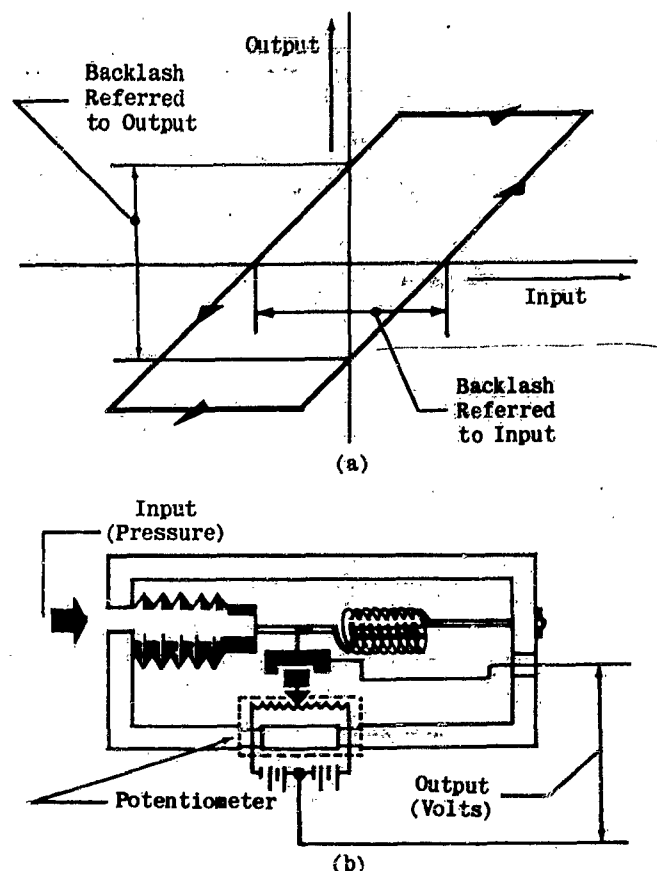


Figure II-4. Backlash Type Discontinuous Non-Linearity

Threshold or flat spot, figure II-6b, occurs in a system equivalent to that of figure II-4 with the addition of a centering spring from the potentiometer wiper to the case. Whenever the bellows has moved beyond the threshold region, the spring holds the wiper against the wiper attachment. Thus when the bellows reverses its direction of travel, there is no lost motion until the threshold region is reached again.

Preload, figure II-6c, is characterized by a step function force or torque versus a displacement away from neutral. Note that the applied force on a device with this characteristic can vary from minus the preload

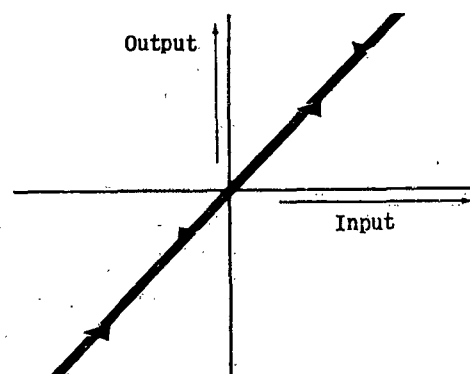


Figure II-5. Backlash

Chapter II
Section 2

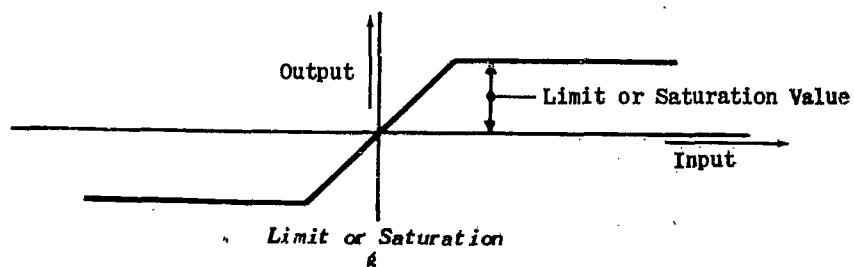
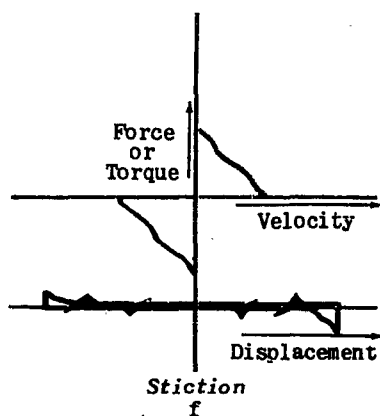
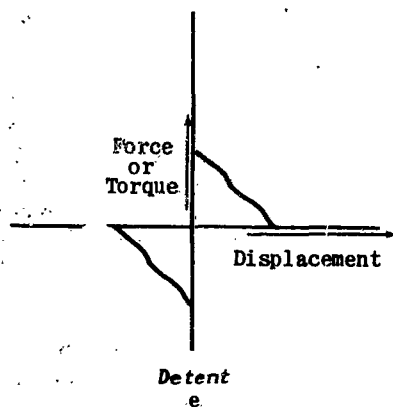
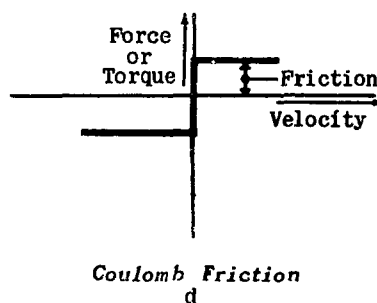
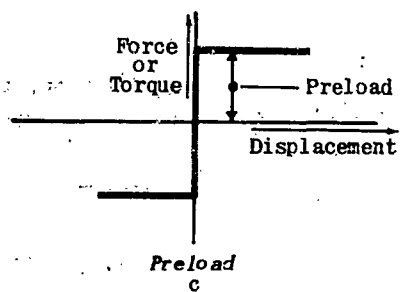
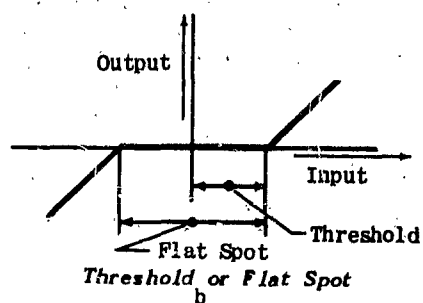
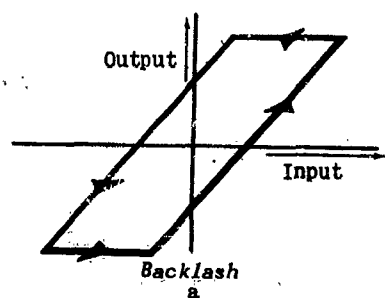


Figure II-6. Discontinuous Non-Linearities

value to plus, without causing any displacement.

Coulomb friction has the same relation to velocity that preload has to displacement, figure II-6d. As shown in the illustration, a force equal to the friction is required to produce a displacement away from zero. To reverse the direction, the force must be reduced to zero, then applied again in equal strength in the opposite direction.

Detent, figure II-6e, differs from preload in that the force reduces to zero at some small value of displacement. The shape of the detent characteristic is determined by design.

Stiction, figure II-6f, is usually assumed to be related to velocity as detent is related to displacement. In general, this assumption is valid only for very low values of velocity. Stiction differs from coulomb friction in that the force reduces to zero at some small value of velocity.

All of the non-linearities so far discussed have concerned behavior of a system near neutral or the null point. An important non-linearity encountered with larger values of input and output is limiting, figure II-6g. This characteristic occurs in any real system. The usual assumption is that the operating range of significance in analysis is "below saturation."

There are many other types of non-linearities that are important in certain control system analyses. But these are discussed in some detail in chapter VI. The essential point to be gathered here is that these examples are typical of what is not acceptable if the system is to be represented by linear differential equations with constant coefficients (II-2).

So far only those effects of coefficients varying with the dependent variable (x) have been considered. The coefficients may also vary with time. In the simple mass-damper-spring example shown above, the temperature of the environment may be changing in some fashion. If this were the case, M , B , and k would also vary with time; then (II-1) would take on the form:

(II-4)

$$f_0(t) \frac{d^n x}{dt^n} + f_1(t) \frac{d^{n-1} x}{dt^{n-1}} + \dots + f_{n-1}(t) \frac{d^{n-1} x}{dt^{n-1}} + \dots + f_n(t) \frac{dx}{dt} + f_{n+1}(t)x = Q(t)$$

This is known as a linear differential equation with variable coefficients. It is more readily solved than a non-linear equation, but the process is much more involved than that required for the simple linear equation with constant coefficients. In particular, it is not amenable to any of the methods to be used in chapters III through V. Consequently, the remainder of this chapter is devoted to mathematical models that can be derived from equations typified by (II-2).

SECTION 3—BLOCK DIAGRAMS AND TRANSFER FUNCTIONS

(a) THE BLOCK DIAGRAM

Engineers have developed methods of working with drawings and diagrams especially designed to provide useful information and to aid in the visualization of certain aspects of a problem. For example, figure II-7 shows a pump geared to a motor whose field voltage is supplied by a remotely located control box. This diagram provides information regarding the number of units composing the system and their relative locations and sizes. That is, it conveys a description of some of the external features of the system. However, it does not provide a thorough understanding of how the system operates. The operation may be seen more

readily in the schematic diagram of figure II-8.

This figure shows an amplifier supplying a voltage, V_a , to the control field of a motor. This voltage alters the motor torque, thus effecting a change in motor speed, n_1 , with a corresponding perturbation of the rate of flow of fluid, Q_0 , in the outlet pipe. The flow of fluid impinges on the flowmeter vane which is balanced by a spring. A potentiometer attached to the vane puts out a voltage, V_p , proportional to vane deflection (and consequently, proportional to Q_0). The flowmeter potentiometer is connected to the control potentiometer in a bridge circuit. Thus, when the control potentiometer is turned clockwise the voltage,

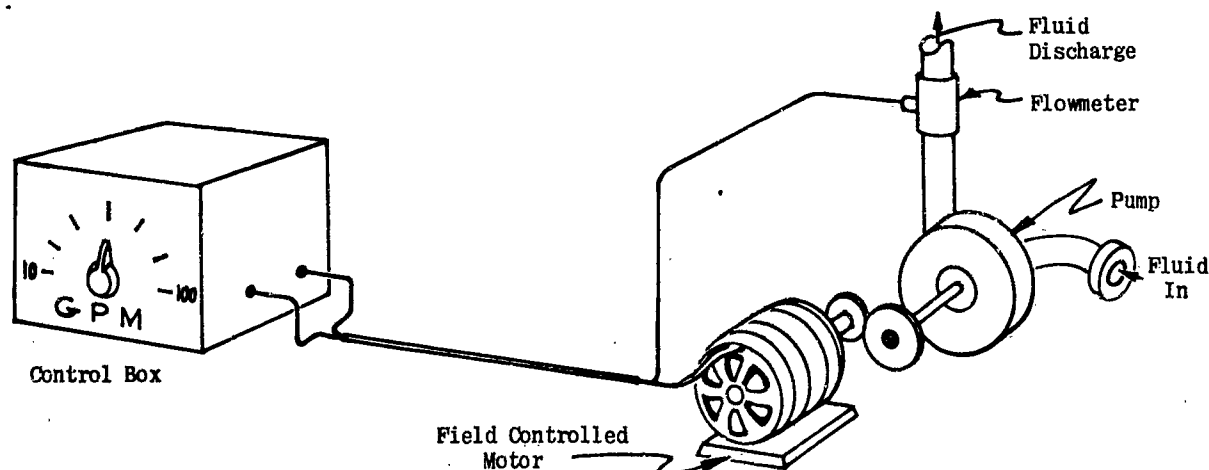


Figure II-7. Pump Drive System

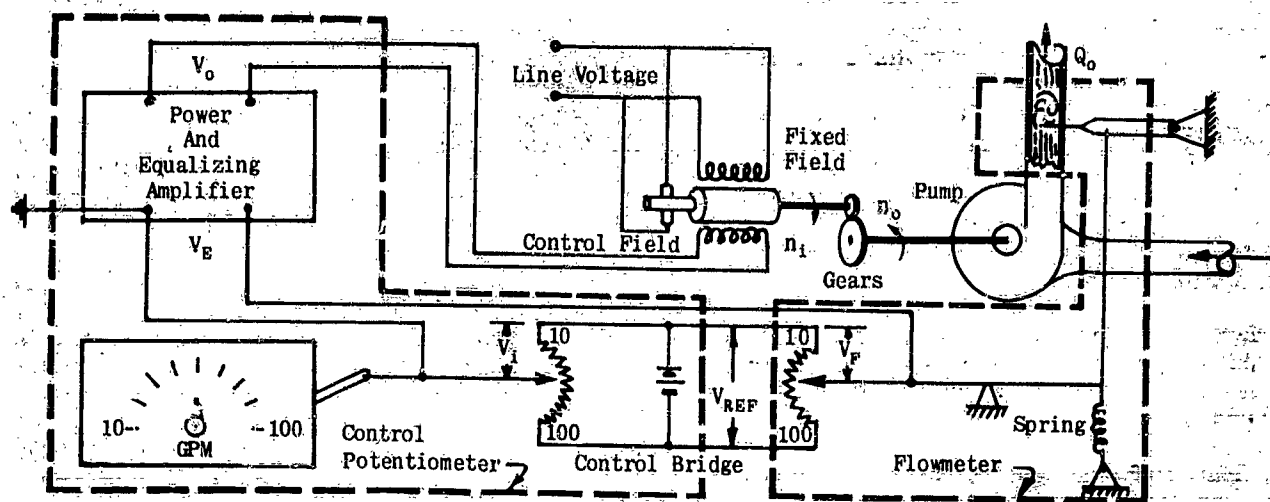


Figure II-8. Schematic Diagram of Pump Drive System

V_i , is increased so that a voltage $V_E = V_i - V_f$ appears at the amplifier input. The amplifier contains special circuitry which increases the current in the motor control field until the flowmeter voltage, V_f , just balances the control potentiometer voltage, V_i , (i.e., $V_i - V_f = 0$). In this way precise control is maintained over the flow rate Q_o .

Although figure II-8 is an abstraction that bears only a slight resemblance to figure II-7, it is useful because it shows the functional relationships of the entire system. Thus, it serves as an aid to the formulation of the mathematical model of the system. In fact, one of the equations has already been derived, namely, $V_E = V_i - V_f$. This equation states that the amplifier input voltage is the difference between the control potentiometer setting and the vane potentiometer position. This relation will be used later in deriving the complete equation of the system.

An even more abstract representation of the system can be devised which allows direct determination of the mathematical model. Consider first the amplifier. It has as its input, V_E , the difference between the control potentiometer voltage and the flow meter potentiometer voltage. Its output is the motor control field voltage, V_o . This can be represented as in figure II-9.

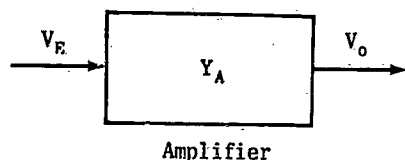


Figure II-9. Block Representation of Amplifier

The term Y_A in the figure is referred to as the amplifier transfer function because it transfers the input to the output, and it is a mathematical expression of the amplifier performance. That is, the output of the unit is given by $V_o = Y_A V_E$.

The motor operation may be expressed in a similar fashion by a block as in figure II-10.

Notice here that the motor input voltage is the amplifier output voltage and that the motor output is the rotational

speed n_1 . The function Y_M transfers the input V_o to the output n_1 .

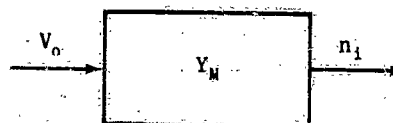


Figure II-10. Block Representation of Motor

There is another input to both the motor and the amplifier, namely, the line voltage. See figure II-8. However, this voltage is considered as merely an energy source of constant magnitude. If this assumption is true, only the steady state characteristics of the system are affected. In this particular case, line voltage is important because it determines the steady state speed of the motor. The question of primary importance in the analysis treated in this book is: "What happens to the system when it is disturbed slightly from the steady state operating point?" Consequently, the transfer function is concerned only with those inputs which will help answer this question. An important assumption upon which all the analysis treated in this volume is based is that a static analysis has previously been made and it has been ascertained that the system is capable of operation at the steady state operating point under consideration.

Returning to figure II-8, the gear unit and the pump may also be represented by blocks, each with a transfer function relating the output to the input, figure II-11.

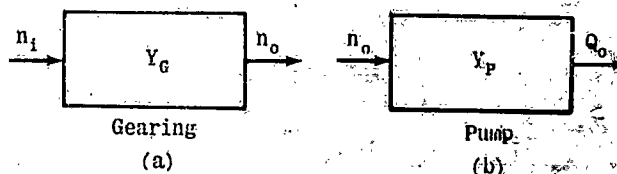


Figure II-11. Block Representations

In figure II-11, the output quantities are defined by $n_o = Y_G n_1$ and $Q_o = Y_P n_o$.

Combining figures II-9, II-10, and II-11 results in figure II-12:

This combination of blocks is known as a block diagram. It shows the cascade arrangement of the amplifier, motor, gearing, and pump. As yet, this block diagram is not a complete functional representation of the pump drive system since the function of the flowmeter in the system is not yet accounted for. However, the diagram does show the transfer of an input (control) quantity, V_E , through various functional components into an output (controlled quantity) fluid flow, Q_o , and on this basis could be called a control system. But regarded

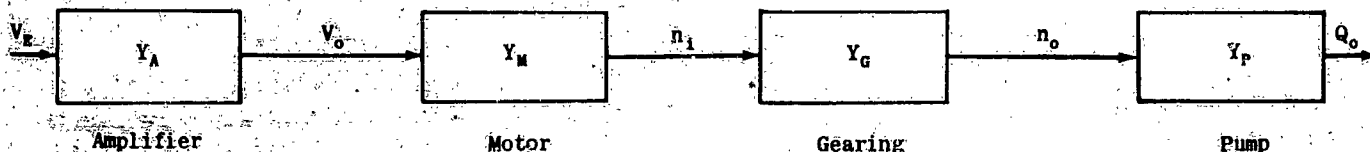


Figure II-12. Block Diagram of Portion of Pump Drive System

as a system, it has the serious deficiency that any variations of the pressure of the fluid supply to the pump produces a change in the output quantity, Q_o . It is evident that to maintain a constant output flow, Q_o , with this system, some means would be needed to vary the input, V_E , whenever changes occur in Q_o .

This means is provided by a flowmeter used as shown in figure II-8. The flowmeter, driven by the fluid, turns a potentiometer, which produces a voltage, V_F . This is represented by figure II-13.

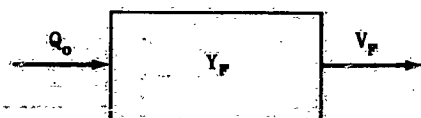


Figure II-13. Block Representation of Flowmeter

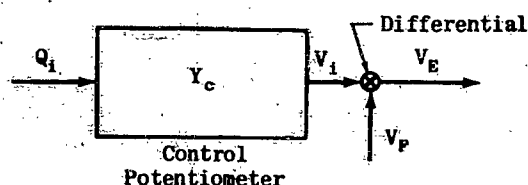


Figure II-14. Block Representation of Input to Amplifier

The position of the control potentiometer wiper, and thus its voltage, is proportional to the desired flow, $V_i = Y_C Q_i$. By virtue of this circuitry, the voltage into

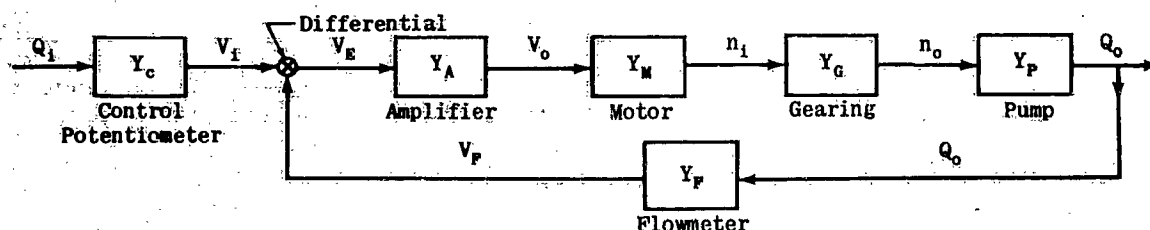


Figure II-15. Block Diagram of Pump Drive System

the amplifier is the difference between the flowmeter potentiometer output, V_F , and the control potentiometer setting, V_i . This is indicated by the scheme shown in figure II-14. The illustration also introduces a new

symbol \oplus , called a differential, which represents the equation $V_E = V_i - V_F$.

The complete block diagram for the pump drive system is derived by combining figures II-12, II-13, and II-14, and is shown in figure II-15. Notice that by following the arrows away from the differential to the right and back through the flowmeter, a complete "loop" is described. Systems that can be represented in this way are known as closed-loop systems and are dis-

tinctly different in their properties from those that can be represented by simple cascading of elements as in figure II-12.

The blocks starting with the amplifier and ending with the pump form what is known as the forward loop elements of the system.* That is, these units transfer the input forward to the output. The flowmeter, which feeds information from the output back to the input, makes up the feedback loop element of the system.

If, in this closed loop system, the output, Q_o , increases for some reason, the feedback voltage, V_F , will also increase. And since the input to the amplifier is $V_E = V_i - V_F$, the amplifier will reduce the motor control field voltage and, consequently, the pump output, Q_o . Automatic regulation of the flow is thus achieved.

If the flowmeter is removed, the feedback from output to input is removed which opens the loop of the block diagram. The system then assumes the form of figure II-12, and is called an "open-loop" system.

A homely example is presented now to emphasize the difference between open and closed loop systems. Consider the automatic washing machine. This device cannot sense the degree of dirtiness of the clothes and performs each operation of its cycle only for a predetermined length of time. This is an open-loop system. On the other hand, when the clothes are laundered by hand, the time and energy expended are a function of the dirtiness of the clothes because the washerwoman

* The terms forward loop, and feedback loop are misleading since "loops" as such are not referred to. However, because of long usage, they are commonly accepted. In this volume, they are also referred to as forward path and feedback path.

Chapter II

Section 3

continuously observes the clothes and controls her behavior according to her degree of satisfaction with them. Thus, she behaves not only according to the command, "Wash these clothes," but also according to what happens as a result of the washing operation. The washerwoman and associated equipment make up a closed loop system.

Whether it refers to an open or a closed loop, the block diagram possesses several important characteristics that deserve special mention. It must be remembered that a block diagram is a functional representation of a physical system. Since the blocks represent functional components rather than physical components, a block may represent several physical units lumped together, or one physical unit may be subdivided into several functional blocks. This combination and separation of physical components into functional units is based upon the operational relationships between the units.

For instance, figure II-15 shows individual blocks for the control potentiometer and amplifier and a symbol for a differential. Actually, the potentiometer and amplifier are located within one control box, while the differential merely represents a method of wiring.

Also in figure II-15, the block representing the motor accounts for the relationship between shaft speed and control field voltage. But since the pump characteristics are important factors in determining shaft speed, certain pump features are actually a part of the motor block. In this example, the flow, Q_0 , is a simple function of pump speed. It is this simple function that is represented by Y_P . It is evident that an intimate knowledge of the behavior of the elements of a system is needed before a block diagram can be constructed. This aspect of the problem will be treated in some detail later.

(b) BLOCK DIAGRAM ALGEBRA

It is a comparatively simple job to derive the mathematical models (transfer functions) of the individual elements such as Y_A , Y_M , Y_G of figure II-15. Each model describes the behavior of the corresponding individual element. To determine how these elements perform when linked together requires the derivation of a new mathematical model that describes the complete closed loop system of figure II-15. This new model is more complex than any represented by the individual elementary transfer functions, and must be examined to obtain a description of the behavior of the entire system. The field of controls system analysis consists simply of a number of special ways of investigating the properties of the closed loop. All of these special methods depend upon manipulating block diagrams such as figure II-15 into simple form. Such manipulations fall into the subject of block diagram algebra.

As systems become more complex, they become unwieldy mathematically. This situation is remedied

by a system of block diagram algebra that makes it possible to reduce even the most complex system to a single block.

This is illustrated by the pump control system discussed above.

The transfer functions Y_A , Y_M , Y_G and Y_P are the transfer functions of the forward elements of the system. Referring to figure II-12, $Y_A = V_o/V_E$, $Y_M = n_1/V_o$, $Y_G = n_o/n_1$, $Y_P = Q_o/n_o$; and since

$$Q_o/V_E = (V_o/V_E)(n_1/V_o)(n_o/n_1)(Q_o/n_o) = Y_A Y_M Y_G Y_P = Y_{FL}$$

the four blocks of the forward path may be reduced to a single block, defined by a forward transfer function Y_{FL} . Figure II-16 - Figure II-15 can then be represented as in figure II-17.



Figure II-16. Equivalent Single Block Representing Forward Path of Figure II-15

Note that the feedback element of figure II-17 is unchanged from that of figure II-15, because the flowmeter transfer function, Y_F , is the only element in the feedback link; consequently it defines the feedback transfer function of the system.

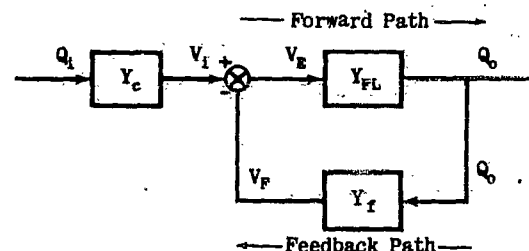


Figure II-17. Simplified Block Diagram of Pump Drive System

The overall relationship between the output, Q_o , and input, Q_i , is obtained in terms of the forward and feedback transfer functions as follows:

The voltage to the amplifier was given by (II-5)

$$V_E = V_i - V_F$$

and since $V_E = Q_o/Y_{FL}$, $V_i = Y_c Q_i$, and $V_F = Y_F Q_o$, these expressions can be substituted into (II-5) yielding $Q_o/Y_{FL} = Y_c Q_i - Y_F Q_o$. Rearranging the terms, $Q_o(1/Y_{FL} + Y_F) = Y_c Q_i$. Clearing fractions and forming the ratio of output to input, the equation is

$$(II-6) \quad \frac{Q_o}{Q_i} = \frac{Y_c Y_{FL}}{1 + Y_F Y_{FL}}$$

(II-5) and (II-6) are the two fundamental relationships on which all closed-loop control systems theory is based. (II-5) is called the actuating error equation. It defines the "actuating error" (as represented by V_E) as the difference between the desired quantity (as represented by V_i) and the output quantity (as represented by V_F). By direct substitution of the input and output, and by application of the forward and feed-

Chapter II
Section 3

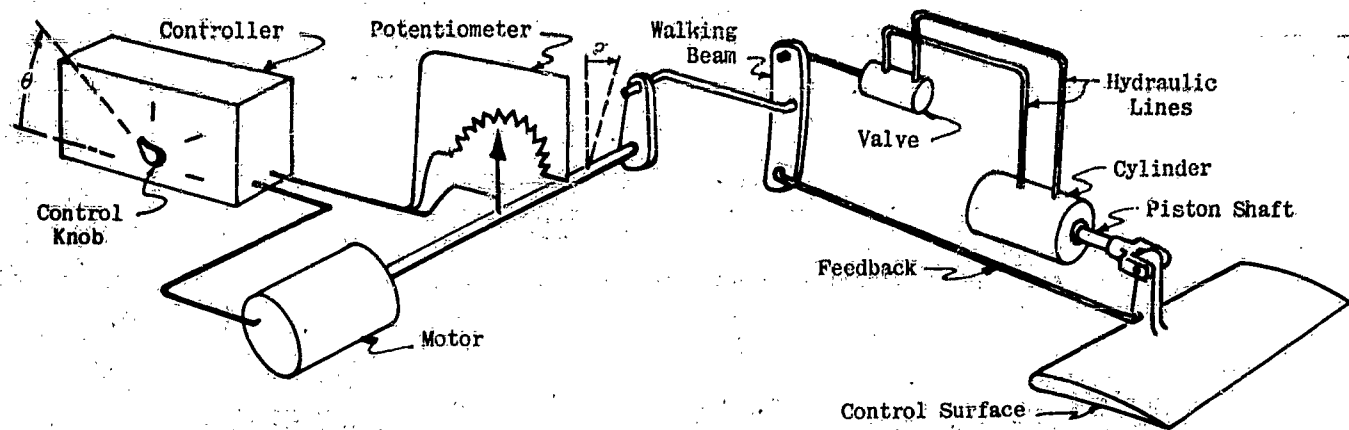


Figure II-22. Surface Control Positioning System

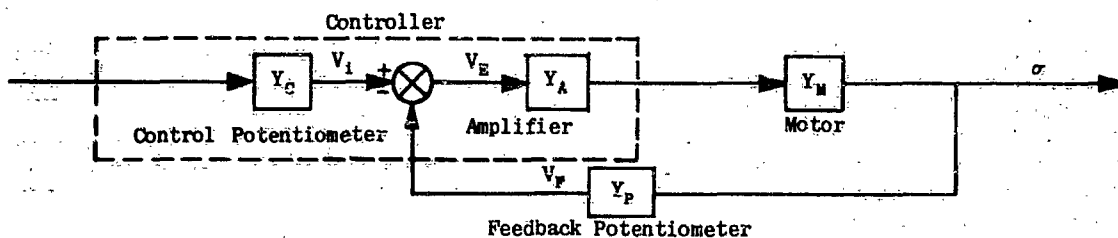


Figure II-23. Motor Control System

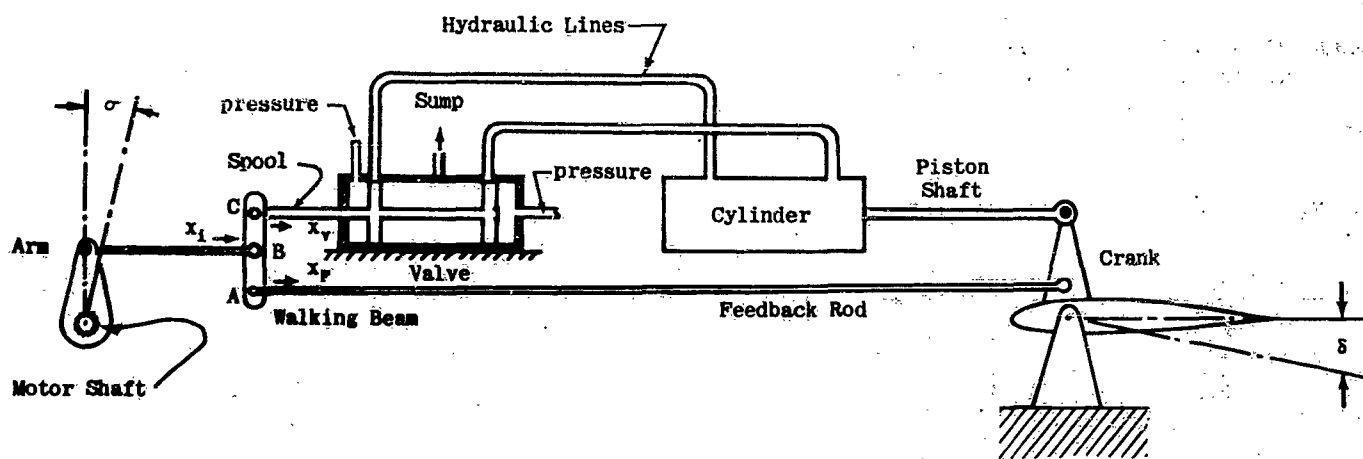


Figure II-24. Hydraulic System

by the block diagram of figure II-25.

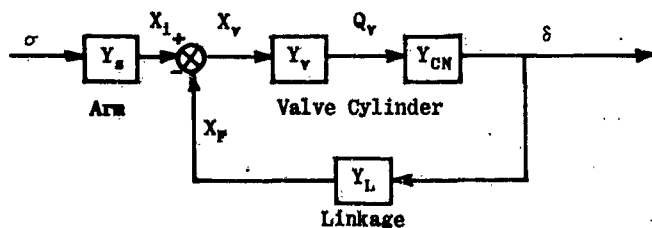


Figure II-25. Hydraulic Control System

The block Y_s is needed to account for the conversion of rotary motion (σ) to linear motion (x_1). The block Y_L accounts for the conversions of surface deflection (δ) to linear feedback displacement (x_f). Included

in both Y_s and Y_L are conversion factors to account for the mechanical advantages of the walking beam.

The block diagram equation of the hydraulic system is

$$(II-10) \quad \frac{\delta}{\sigma} = Y_s \frac{Y_v Y_{CN}}{1 + Y_v Y_{CN} Y_L} \quad \text{or} \quad \frac{\delta}{\sigma} = Y_s Y_H \quad \text{where} \quad Y_H = \frac{Y_v Y_{CN}}{1 + Y_v Y_{CN} Y_L}$$

The complete block diagram of figure II-22 is constructed by combining figures II-23 and II-25. This is shown in figure II-26.

Making use of the definitions of (II-9) and (II-10), the system is reduced to a simple cascade, figure II-27.

By defining $Y_{CP} = Y_C Y_F Y_s Y_H$ the system is reduced to a single block, figure II-28.

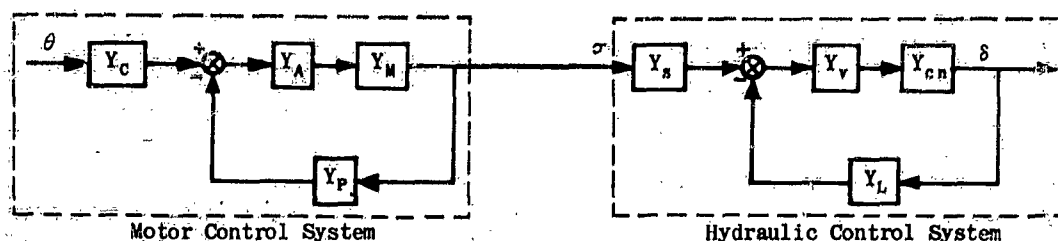


Figure II-26. Surface Control Positioning System



Figure II-27. Surface Control Positioning System

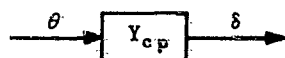


Figure II-28. Surface Control System

A block diagram form of great value in the analysis to be treated in later chapters is shown in figure II-29.

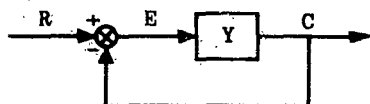


Figure II-29. Unity Feedback Loop

The block diagram equation for this system is derived using the principles of figure II-18. Referring to figure II-29, $E = R - C$ and $C = YE$, therefore $C/Y = R - C$.

$$(II-11) \quad \frac{C}{R} = \frac{Y}{1 + Y}$$

Equation (II-11) should be compared to (II-7). There is one essential difference between these two forms: the presence of a product of two different functions in the denominator of (II-7). Since (II-11) can be derived from (II-7) if the feedback transfer function is unity ($Y_2 = 1$), the system in figure II-29 is referred to as a unity feedback system. For this type of system the closed-loop transfer function is simply

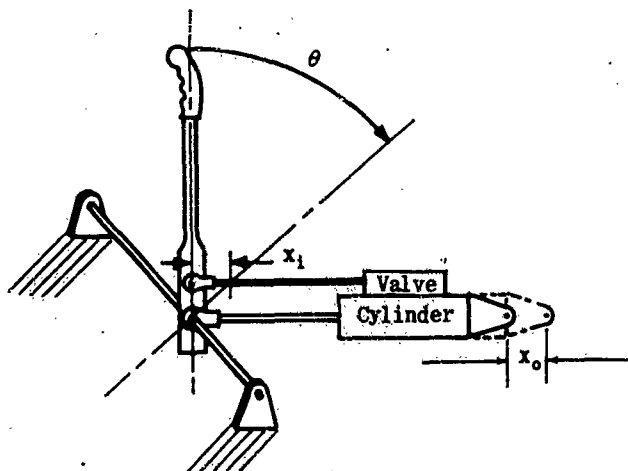


Figure II-30. Unity Feedback Control System

$$(II-12) \quad \text{CLOSED-LOOP TRANSFER FUNCTION} =$$

$$\frac{\text{FORWARD TRANSFER FUNCTION}}{1 + \text{FORWARD TRANSFER FUNCTION}}$$

A common example of an inherently unity feedback type system is the hydraulic system shown in figure II-30.

Since the valve is attached directly to the cylinder, there is a one to one follow-up. The sequence is as follows: control stick rotation displaces the valve spool relative to the valve allowing fluid to flow into the cylinder. The fluid flow is accompanied by cylinder displacement. As the cylinder moves, it carries the valve housing with it. Since the spool is attached to the stick (now held stationary), the cylinder will come to rest when it has carried the valve housing to the neutral position relative to the spool.

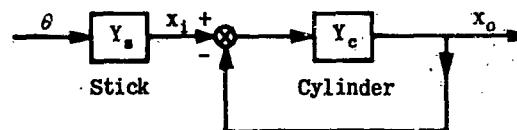


Figure II-31. Unity Feedback Hydraulic System

A control system with non-unity feedback can be converted to the form of figure II-29 by the following sequence. For example, the system appears as in figure II-32; its transfer function may be written as

$$(II-13) \quad \frac{C}{R} = \frac{1}{Y_2} \frac{Y_1 Y_2}{1 + Y_1 Y_2}$$

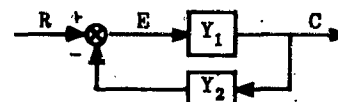


Figure II-32. Non-Unity Feedback System

Replacing $Y_1 Y_2$ by Y , the closed loop of figure II-33 is obtained.

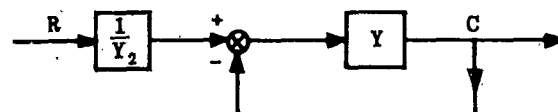


Figure II-33. Unity Feedback System

After the closed loop is analyzed, it is necessary to multiply the results by $1/Y_2$ in order to describe the complete closed loop behavior, i. e. $C/R = (1/Y_2) [Y/(1 + Y)]$.

The unity feedback system is very seldom encountered

in physical networks, but it is very helpful in analyzing closed loop systems.

Another useful diagram is that of the unity forward path systems, figure II-34.

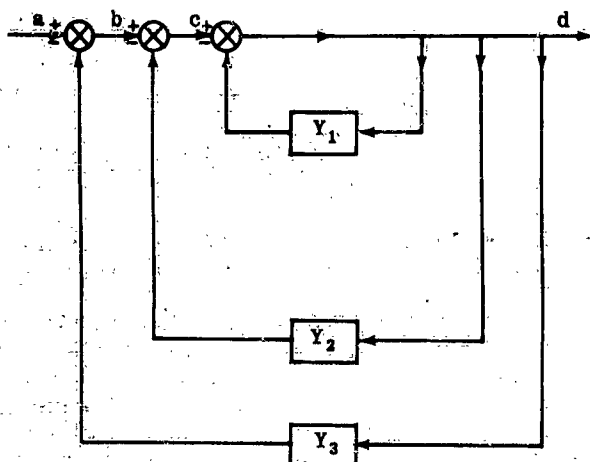


Figure II-34. Unity Forward Path

The transfer function for this system is derived by eliminating one loop at a time. Since the forward transfer function is unity, the first loop is eliminated by applying (II-7) to figure II-34, i.e., $d/c = 1/(1+Y_1)$. The resulting diagram is shown in figure II-35.

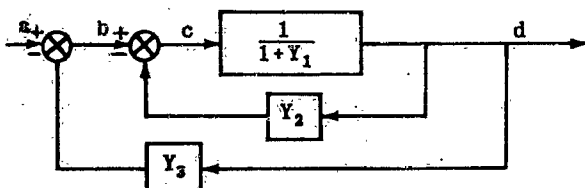


Figure II-35. First Step in Reduction of Figure II-34

Applying (II-7) to the second loop,

$$\frac{d}{b} = \frac{\frac{1}{1+Y_1}}{1+Y_2\left(\frac{1}{1+Y_1}\right)} = \frac{1}{1+Y_1+Y_2}$$

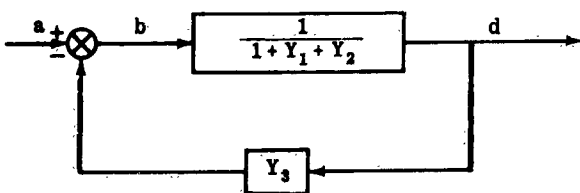


Figure II-36. Elimination of Second Loop of Figure II-34

Finally,

$$\frac{d}{a} = \frac{\frac{1}{1+Y_1+Y_2}}{1+Y_3\left(\frac{1}{1+Y_1+Y_2}\right)} = \frac{1}{1+Y_1+Y_2+Y_3}$$

and the block diagram is shown in figure II-37.

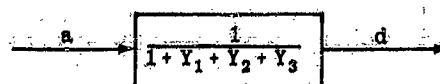


Figure II-37. Final Form of Figure II-34

A symbol useful in developing block diagrams is the adder (\oplus). Application of this symbol is rather subtle since it is mechanized by the same type of apparatus as the differential (\otimes).^{*} One important application is in describing a so called open loop-closed loop control system figure II-38. Here the input is compared to a feedback quantity as in any of the closed loop systems described previously. In addition it is applied directly to the final control element (Y_P) through an "open loop controller" (Y_{OL}). Here the adder symbol is used to show that signals are being added into the loop from an external source (Y_{OL}).

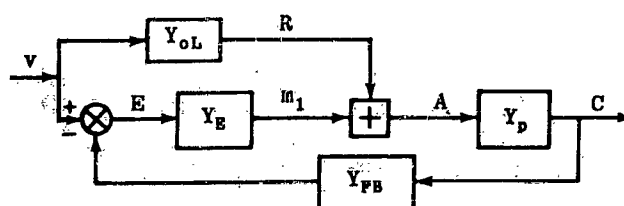


Figure II-38. Open Loop - Closed Loop Controller

The adder is also used to describe positive feedback systems. If the quantity fed back is added to the input, the block diagram will appear as in figure II-39. The closed loop equation now has the form:

(II-14)

$$\frac{C}{V} = \frac{Y_1}{1-Y_1Y_2}$$

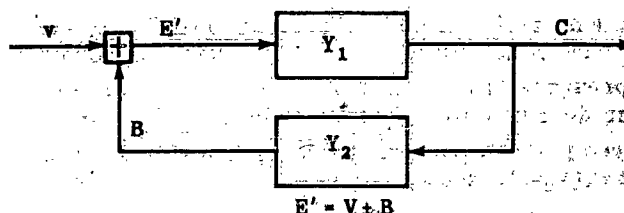


Figure II-39. Positive Feedback System

A third application of the adder is shown in figure II-40. This is a parallel arrangement of blocks. The adder output is given by $d = a' + b' + c'$, and since $a' = aY_a$, $b' = bY_b$, and $c' = cY_c$, $d = aY_a + bY_b + cY_c$. Note that

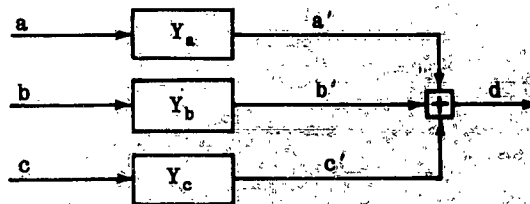


Figure II-40. Parallel Arrangement of Blocks

* See Lauer, Lesnick and Matson (Ref. 3) or James, Nichols and Phillips, (Ref. 4) for excellent descriptive material on mechanization of these functions.

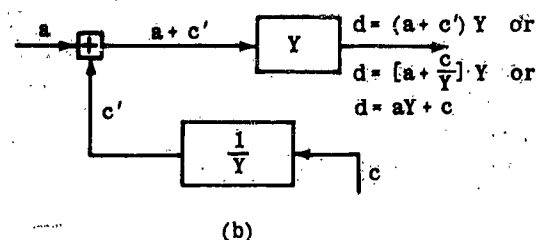
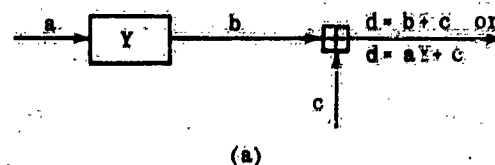
the cascade or series arrangement of blocks, figure II-12, represented a multiplication operation, and the parallel arrangement represents an adding operation.

The location of the adder or differential can be shifted in a block diagram to aid in the simplification process. Figure II-41 illustrates this.

Table II-1* is an extensive list of similar transformation pairs to be used in modifying and reducing complex diagrams. Note that item 10 was used to transform the cascaded blocks of figure II-12 into that of the single block of figure II-16. Also note that figure II-29 was modified to figure II-20 by using item 16.

To illustrate the use of this table a complex multi-loop block diagram, figure II-42, is reduced to a simple single loop as follows:

* This table was adapted from Reference 6



**Figure II-41. Equivalent Block Diagrams—
Moving of Summing Point**

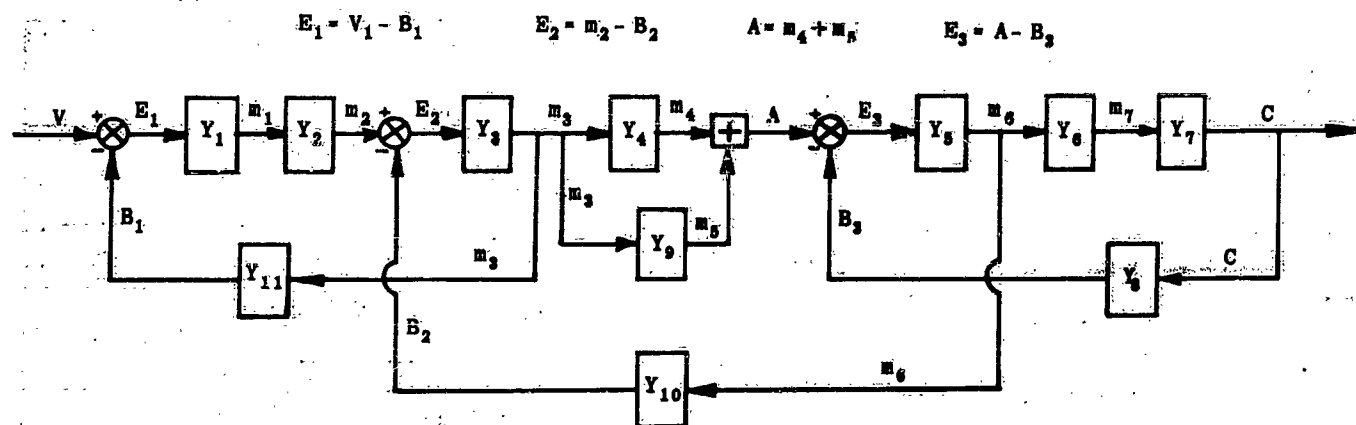


Figure II-42. Multi-Loop Block Diagram

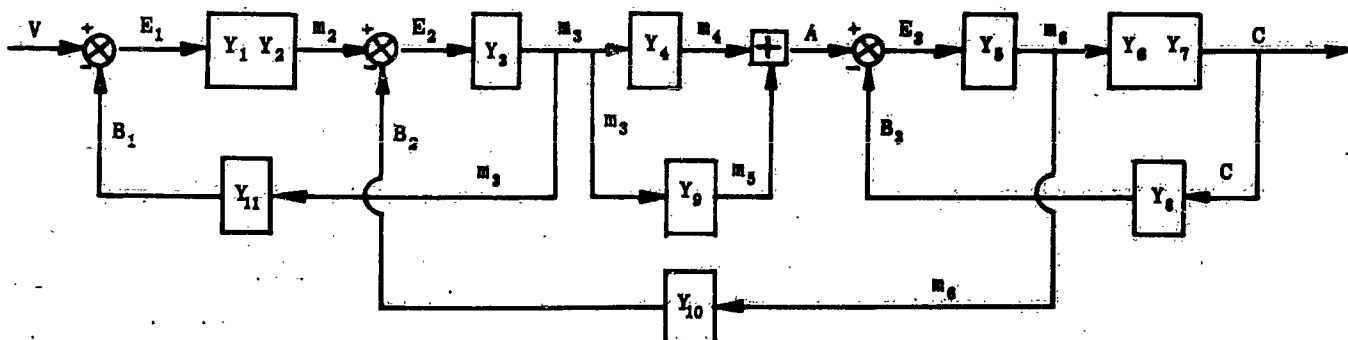
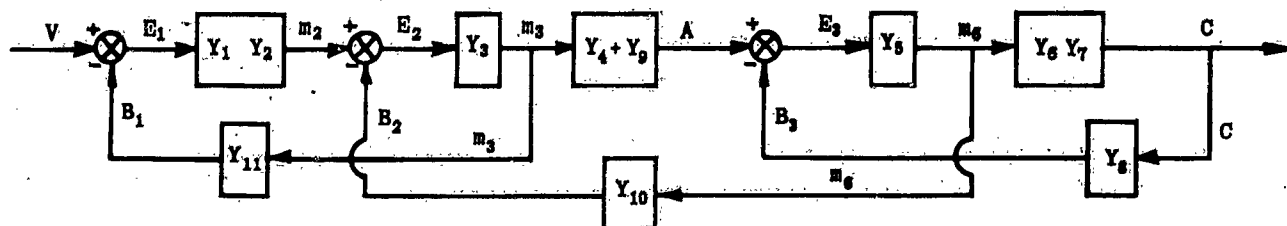


Figure II-43. Cascade Elements of Figure II-42 Combined



**Figure II-44. Forward Loop of Figure II-43
Eliminated**

Chapter II
Section 3

Transformation	Original Diagram	Equivalent Diagram	Equation
1. Interchange of Blocks			$b = aY_1Y_2$
2. Interchange of summing points			$d = a - b + c$
3. Rearrangement of summing points.			$d = a - b - c$
4. Moving a summing point ahead of an element			$d = aY - c$
5. Moving a summing point beyond an element			$c = (a - b)Y$
6. Moving a takeoff point ahead of an element			$b = aY$
7. Moving a takeoff point beyond an element			$b = aY$ $a = b/Y$ $= a$
8. Moving a takeoff point ahead of a summing point			$c = a - b$
9. Moving a takeoff point beyond a summing point			$c = a - b$ $a = c + b$
10. Combining cascade elements			$b = aY_1Y_2$
11. Removing an element from a forward loop			$d = a(Y_1 - Y_2)$
12. Inserting an element in a forward loop			$d = aY_1 - a$
13. Eliminating a forward loop			$d = a(Y_1 - Y_2)$
14. Removing an element from a feedback loop			$d = \frac{aY_1}{1 + Y_1Y_2}$
15. Inserting an element in a feedback loop			$d = \frac{aY_1}{1 + Y_1}$
16. Eliminating a feedback loop			$d = \frac{aY_1}{1 + Y_1Y_2}$
16a.			$d = a \frac{Y_1}{1 + Y_1}$
16b.			$d = a \frac{1}{1 + Y_1}$
17. Inserting a feedback loop			$d = aY_1$
17a.			$d = aY_1$

Table II-1

As a starting point, the series blocks are combined wherever possible (item 10). Thus Y_1 and Y_2 , and Y_6 and Y_7 are combined. Figure II-43 shows the first reduction.

Next, the forward loop between m_3 and A is simplified (transformation 13; see figure II-44).

Figure II-45 shows the relocation of the take-off point and differential point of the main (outer) loop (transformations 7 and 4). At the original take-off point figure II-44, the input to the feedback block denoted by Y_{10} was m_6 . The take-off point is now at C. Consequently, C must be transferred back to m_6 through the block $1/Y_6Y_7$ before being fed into Y_{10} . Similarly, B_2 , the output of Y_{10} , must be transferred to $B'_2 = B_2/Y_1Y_2$

before being subtracted from the input V. Thus, the input to Y_1 is $(V - B'_2 - B_1)Y_1Y_2$. Substituting for B'_2 , the input is $(V - B_2/Y_1Y_2 - B_1)Y_1Y_2 = (V - B_1)Y_1Y_2 - B_2$. Since $(V - B_1)Y_1Y_2 = m_2$ it can be seen that this is the same input to Y_3 as in figure II-44.

The three sets of cascaded blocks are combined in figure II-46, leaving two minor loops separated by a single block, all enclosed by a major loop.

The two inner loops are easily converted into single blocks (transformation 16) as indicated in figure II-47. From this, the final modified diagrams of figure II-48 are evolved.

The final diagram allows the direct determination of

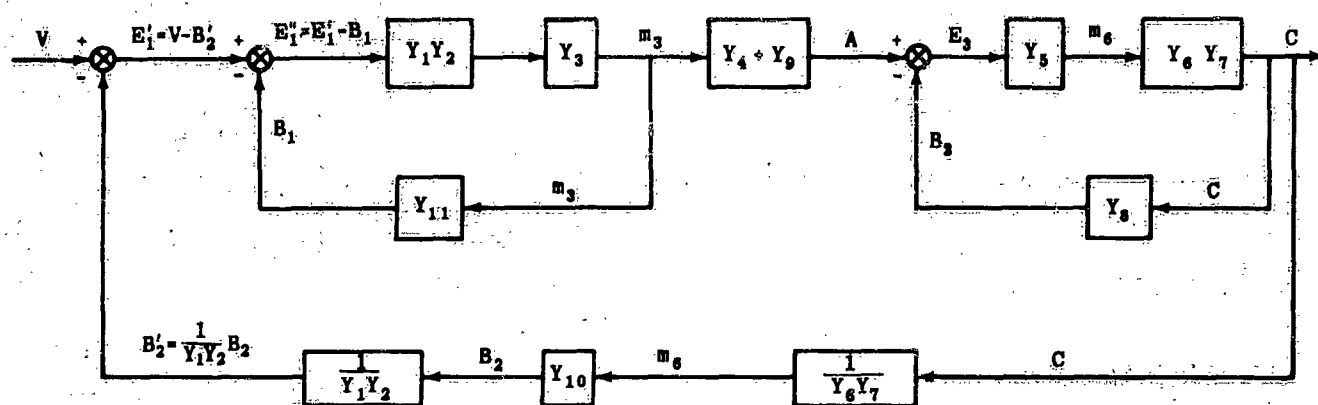


Figure II-45. Take-Off Point for m_6 Moved to Right and Summing Point for B_2 Moved to Left

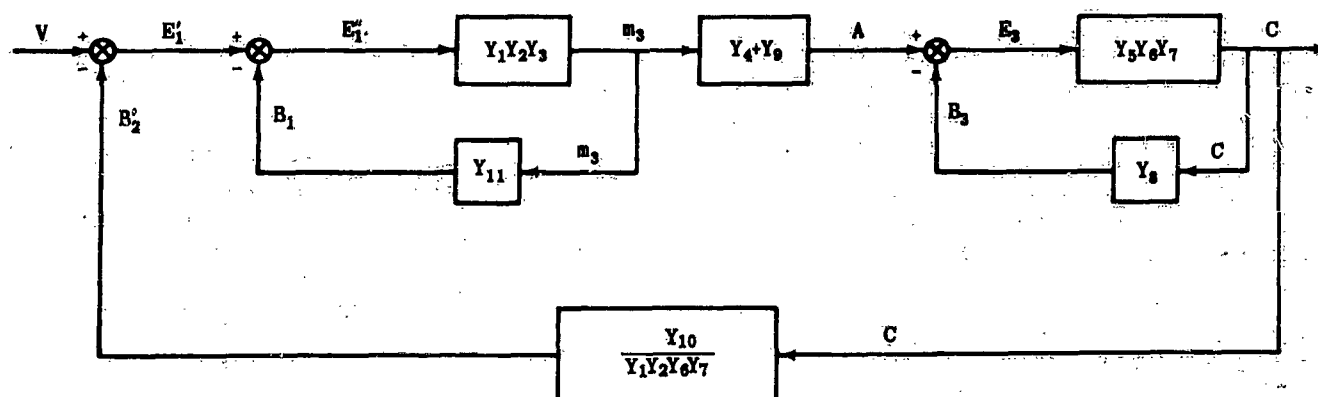


Figure II-46. Cascade Elements of Figure II-45 Combined

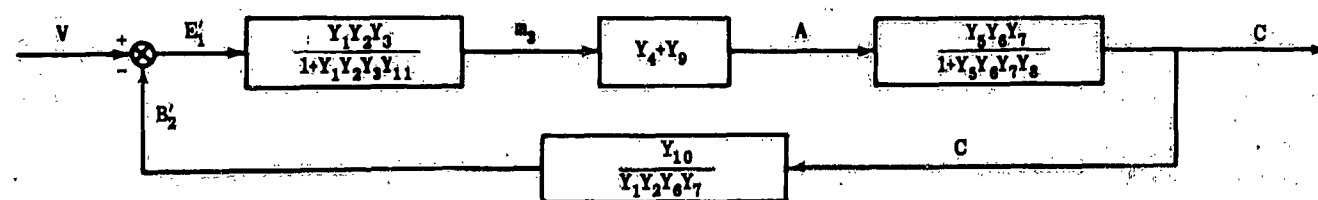


Figure II-47. Two Inner Loops of Figure II-46 Eliminated

Chapter II

Section 3

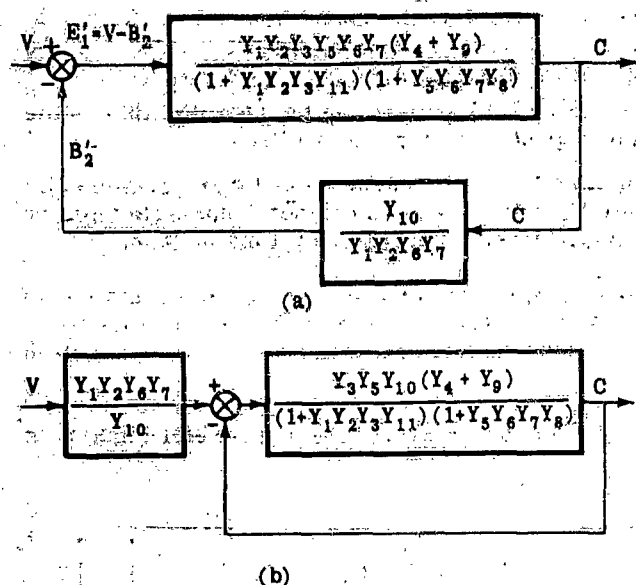


Figure II-48. Final Simplified Block Diagrams of Figure II-42

the overall closed loop equation C/V . For this diagram, it is given by:

$$(II-15) \quad \frac{C}{V} = \frac{Y_1 Y_2 Y_3 Y_5 Y_6 Y_7 (Y_4 + Y_9)}{V (1 + Y_1 Y_2 Y_3 Y_{11}) (1 + Y_5 Y_6 Y_7 Y_8) + Y_3 Y_5 Y_{10} (Y_4 + Y_9)}$$

(c) EQUIVALENT BLOCK DIAGRAM

The transfer function (Y) for many physical units can be derived directly by ordinary methods of analysis. However, it is often advantageous to construct block diagrams from the equations as the derivation proceeds. This not only helps in the derivation by indicating subsequent steps, but also greatly aids in developing an understanding of the physical situation that the equations represent. The procedure is illustrated here by deriving the transfer functions of a compound wound d-c motor with a separately excited control field.

Figure II-49a is a schematic representation of the motor. This can be simplified as in figure II-49b. In figure II-49b resistances and inductances are lumped together in the respective branches of the circuit. By this simplification, analysis is facilitated without loss of specific information. It is important to note that positive control field flux is in opposition to the fixed shunt field.

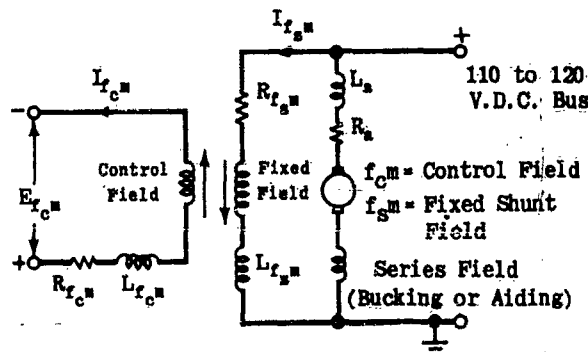
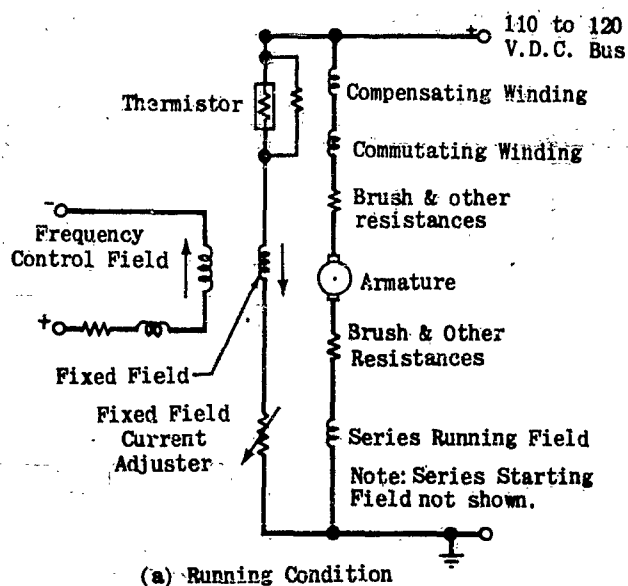
The following assumptions are made for this motor analysis:

- The system is linear.
- Perturbations about operating points are small.
- There exists negligible coupling other than that intended between circuits, components, or voltage and frequency control systems.
- Negligible armature reaction effects.
- Negligible d-c bus voltage variations.
- Negligible shunt field current variations.
- Negligible saturation of magnetic circuits.

The motor constants are defined as follows:

Total armature back EMF = E_a

Total armature current = I_a
 Total d-c bus voltage = $E_{d.c.}$
 Total flux = Φ
 Motor speed = θ
 Armature circuit resistance = R_a
 Armature circuit inductance = L_a
 Control field resistance = $R_{f.c.}$
 Control field inductance = $L_{f.c.}$
 Total "effective" motor control field current = $I_{f.e.}$
 Total motor torque = T_m
 Motor-load inertia = J
 Motor-load viscous friction = B
 Small perturbations are denoted by lower case letters.



Note: Increasing $I_{f.c.}$ Decreases Net Flux

(b) Reduced Diagram of (a).

Figure II-49. Motor Schematic

Then, from motor theory, $E_m = K_1 \Phi \theta$. In addition to being a function of fixed and control field currents, Φ is also a function of armature current since the motor is compounded. Then, perturbations in E_m are denoted by

$$(II-16) \quad e_m = \frac{\partial E_m}{\partial \Phi} \left(\frac{\partial \Phi}{\partial I_a} i_a + \frac{\partial \Phi}{\partial I_{f.c.}} i_{f.c.} \right) + \frac{\partial E_m}{\partial \theta} \theta$$

(No term due to fixed field current exists since its perturbations are assumed negligible.) Here it is important to note that the control field is a bucking

field, that is $(\partial \phi) / (\partial I_{f_c})$ is by convention a negative number at the operating point.

Motor theory also relates armature current to back EMF as follows: $E_{a_c} = E_m + I_a (R_a + sL_a)$ where $s = d/(dt)$ or for small perturbations,

$$(II-17) \quad i_a = -\frac{e_m}{R_a(\tau_a s + 1)}$$

where $\tau_a = L_a/R_a$

Equations (II-16) and (II-17) indicate that the back EMF component is a feedback system as shown in figure II-50.

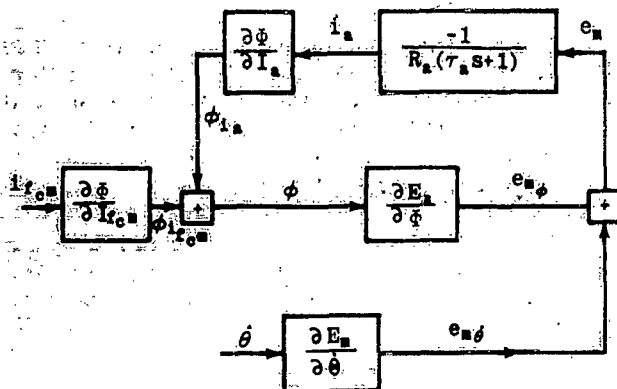


Figure II-50. Back EMF Block Diagram

Motor torque is proportional to the product of the net flux and the armature current, i.e., $T_m = K_2 \phi I_a$.

Because of the series field, perturbations in ϕ are a function of armature current I_a as well as the control field current I_{f_c} . Consequently for small perturbations

$$(II-18) \quad t_m = \frac{\partial T_m}{\partial \phi} \left(\frac{\partial \phi}{\partial I_a} i_a + \frac{\partial \phi}{\partial I_{f_c}} i_{f_c} \right) + \frac{\partial T_m}{\partial I_a} i_a$$

Figure II-51 is a block diagram of the motor torque equation.

The important dynamic loads on the motor are simply motor-load inertia and damping, the latter caused

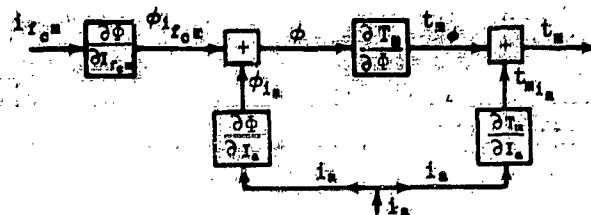


Figure II-51. Motor Torque Block Diagram

chiefly by windage. Therefore, $T_m = Js \dot{\theta} + B \dot{\theta}$.

In perturbed values

$$(II-19) \quad \frac{\dot{\theta}}{t_m} = \frac{1}{B(\tau_B s + 1)}$$

where $\tau_B = J/B$.

Figure II-52 shows the speed-torque block diagram.

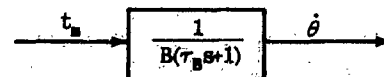


Figure II-52. Speed-Torque Block Diagram

The control field transfer function does not affect motor stability, however it is presented here to complete the description of the motor "components."

Input to the motor is the control field voltage E_{f_c} . Because of the control field inductance L_{f_c} there is a time lag between the control field current I_{f_c} and the control field voltage E_{f_c} .

$$i_{f_c} = \frac{e_{f_c}}{L_{f_c}s + R_{f_c}} = \frac{K_{f_c} e_{f_c}}{\tau_{f_c} s + 1}$$

where $\tau_{f_c} = (L_{f_c}) / (R_{f_c})$ and $K_{f_c} = 1 / (R_{f_c})$. This is shown in figure II-53.

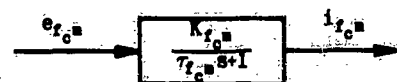


Figure II-53. Motor Control Field Block Diagram

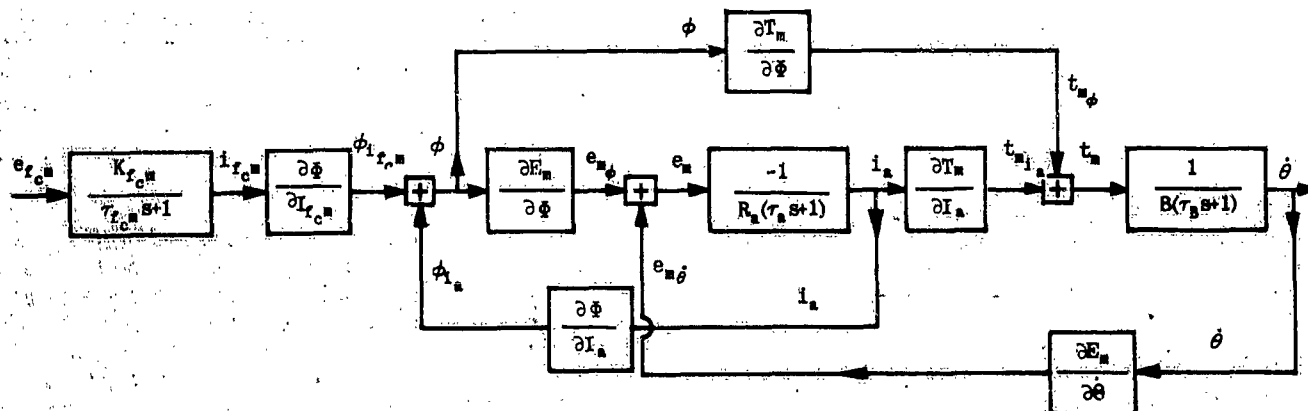


Figure II-54. Motor System Block Diagram

Combining figures II-50, II-51, and II-52, the motor system block diagram of figure II-54 is obtained. The system transfer function is derived by combining equations (II-16), (II-17), (II-18), and (II-19), and is shown by equation (II-20).

$$(II-20) \quad Y_m(s) = \frac{\theta}{E_{rcm}} = \left[\frac{K_{rcm}}{T_{cm}s + 1} \right] \times \frac{\partial \Phi}{\partial I_{rcm}} \left[\frac{\left[\frac{\partial T_m}{\partial \Phi} L_a \right] s + \left[\frac{\partial T_m}{\partial \Phi} \left(R_a + \frac{E_m}{\partial \Phi} \frac{\partial \Phi}{\partial I_a} \right) - \frac{\partial E_m}{\partial \Phi} \left(\frac{\partial T_m}{\partial I_a} + \frac{\partial \Phi}{\partial I_a} \right) \right]}{L_a J \left[s^2 + \left\{ \frac{1}{L_a} \left(\frac{\partial E_m}{\partial \Phi} \frac{\partial \Phi}{\partial I_a} + R_a \right) + \frac{B}{J} \right\} s + \frac{1}{J L_a} \left\{ B \left(\frac{\partial E_m}{\partial \Phi} \frac{\partial \Phi}{\partial I_a} + R_a \right) + \frac{\partial E_m}{\partial \Phi} \left(\frac{\partial T_m}{\partial \Phi} \frac{\partial \Phi}{\partial I_a} + \frac{\partial T_m}{\partial I_a} \right) \right\} \right]} \right]$$

Equation (II-20) is formidable to say the least. However, figure II-54, which is its block diagram form, is fairly straightforward in that it shows in readily identifiable form the manner in which all of the electrical and magnetic entities of the motor interact to produce the dynamic behavior defined in the transfer function equation (II-20).

(d) THE TRANSFER FUNCTION

In subsection (b) of this chapter the terms Y_A , Y_m , and Y_p of the pump drive system were denoted as the block transfer functions of the functional components. In subsection (c), a means of obtaining block diagrams from differential equations was discussed with the derivation of specific transfer functions considered incidental to the process. These sections were primarily concerned with developing block diagram abstractions of physical systems. An ultimate aim of this volume is to provide the designer with the basic tools required to perform experiments on paper. The block diagram abstraction is a necessary step in this direction.

To perform experiments on paper, the system designer must have available mathematical models of system components which

1. completely define the performance of the components,
2. can be combined with other models according to the procedure developed in section II-3b to obtain models of a system.

The transfer function is a form of mathematical model fulfilling the above requirements. This subsection is concerned with the derivation and interpretation of transfer functions. In the process of explanation, it is shown that the transfer functions can be combined and do completely define the performance of a system.

To simplify the presentation, the characteristics of transfer functions are developed by the extensive use of examples. Transient and stability characteristics are discussed for specific systems, with pertinent generalizations noted without proof.

Transfer functions of linear systems with constant parameters* can always be expressed as the ratio of two polynomials in the Laplace transform variable s .** Further, it is always possible to write

* A linear system is one whose properties may be expressed mathematically in terms of linear differential equations.

** The Laplace transform is extensively used in this volume. The reader who is unfamiliar with this method is referred to Reference 5.

$N(s) / D(s) = KG$ in which K is a positive constant and G is a function of s expressed as a non-dimensional ratio of products of the factors of $N(s)$ and $D(s)$. The roots of $N(s)$ and $D(s)$ are referred to as the zeros and poles, respectively, of the transfer function, and to-

gether with the "gain" K they completely define the system.

The simplest transfer function is one in which G is unity. Consider, for example, the potentiometer, gear box, and amplifier of the pump drive system of subsection (a). The potentiometer is calibrated so that a certain setting representing a desired fluid flow rate, Q_1 , will produce a voltage, V_1 . Since the potentiometer is essentially linear, the transfer function is $V_1/Q_1 = K_c$ volts/gpm, where K_c is a constant determined by the reference voltage, V_{ref} , and the calibration of the potentiometer. The potentiometer block is now given by figure II-55.



Figure II-55. Potentiometer Block

The gear transfer function is simply the gear ratio since inertias and friction effects have been lumped into the motor block. Denoting the gear ratio as K_g , the transfer function is then $n_o/n_i = K_g$ (dimensionless).



Figure II-56. Gear Block

Ideally, the transfer function for the amplifier is defined by its gain, K_A . For an input voltage, V_E , the output voltage, V_O , is equal to $K_A V_E$. The transfer function is K_A volts/volt and is constant. It is quite possible, however, that the output voltage will not vary instantaneously as the input voltage varies. That is, there may be an elapsed time, or time lag, between the change in input voltage and the proportional change in the output voltage. In this case, the G would not be unity. A detailed study of the amplifier would result in an analytical expression for G which would account for the time lag.

For a simple example of a component which has a G function other than unity, consider the hydraulic amplifier shown in figure II-57. An input displacement in the direction shown will open the valve so that the fluid will flow in the direction of the arrows. The result is that the load is displaced to the right. From the figure, the following relationship can be derived:*

* For input x_1 , pivot is at A; for output x_o , pivot is at B.

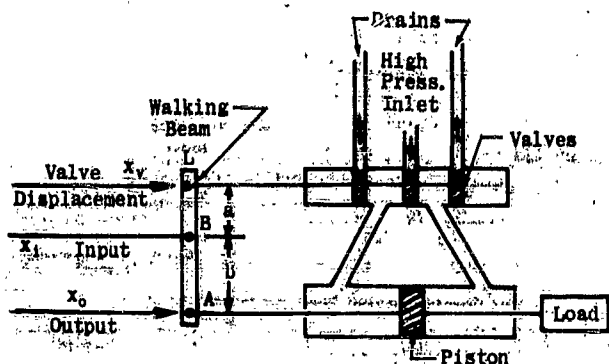


Figure II-57. Schematic of Idealized Hydraulic Amplifier

$$(II-21) \quad x_v = \frac{x_1}{b}(a+b) - x_o \frac{a}{b}$$

The valve displacement x_v controls the quantity of fluid flowing through the inlet tube. This also governs the velocity of the piston displacement. When the hydraulic fluid is incompressible, this relationship is defined by the equation

$$(II-22) \quad \frac{dx_o}{dt} = Cx_v$$

or, transformed, $Cx_v = sx_o$, where s is the Laplace operator and C is the piston velocity per unit valve displacement (constant for constant inlet fluid pressure). These two equations are illustrated in figure II-58. Combining the two figures yields the closed-loop block diagram for the hydraulic amplifier, figure II-59. Note that the walking beam is broken up into two functional parts: the input function and the feedback function.

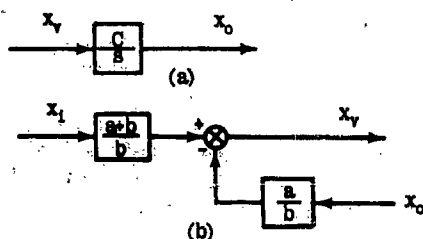


Figure II-58. Components of Hydraulic Amplifier

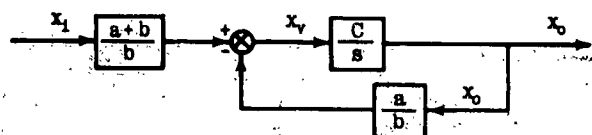


Figure II-59. Closed-Loop Diagram of Hydraulic Amplifier

The output-input transfer function is given by

$$(II-23) \quad \frac{x_o(s)}{x_1(s)} = \frac{a+b}{b} \left(\frac{\frac{C}{s}}{1 + \frac{a}{b} \frac{C}{s}} \right) = \frac{a+b}{a} \left(\frac{1}{\frac{b}{aC}s + 1} \right)$$

$$\frac{x_o(s)}{x_1(s)} = \frac{K_L}{\tau_v s + 1}$$

where $K_L = (a+b)/a$, $\tau_v = b/(aC)$.

Equation (II-23) is illustrated in figure II-60. This diagram is functionally equivalent to figure II-59, but note that the physical relationships are lost in figure II-60.

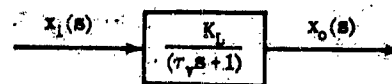


Figure II-60. Block Representation of Hydraulic Amplifier

The transfer function (II-23) represents a single time lag. To illustrate this point, the response of the hydraulic system to an input, x_1 , will be determined.

$$(II-24) \quad x_o(s) = x_1(s) \frac{K_L}{\tau_v} \frac{1}{s + \frac{1}{\tau_v}}$$

A common type of input that is used to evaluate a system is the unit step function (heavy line in figure II-61). That is,

$$(II-25) \quad x_1(t) = \begin{cases} 0 & t < 0 \\ 1 & t > 0 \end{cases}$$

The Laplace transform of (II-25) is $x_1(s) = 1/s$. Consequently, equation (II-24) becomes

$$(II-26) \quad x_o(s) = \frac{K_L}{\tau_v} \frac{1}{s(s + \frac{1}{\tau_v})}$$

The inverse transform of (II-26) is: $x_o(t) = K_L(1 - e^{-t/\tau_v})$. This expression is plotted in figure II-61 for two values of τ_v where $\tau_2 > \tau_1$.

Notice that K_L (the "gain") determines the ratio of the output to the input in the steady state. It is clear from figure II-61 that the quantity τ_v determines how fast the output approaches a steady state value. A practical figure of merit is obtained by letting $t = \tau_v$, then $x_o = 0.633K_L$. So, τ_v represents the time required for the output to reach approximately 63.3% of its steady state value. It is convenient to call τ_v the "time constant" of the element.

Evidently a system whose transfer function is of the form of (II-24) does not respond instantaneously to changes in the input. Since the time delay is represented by a first order equation the system is said to have a "first order time lag."

Another important observation is that the quantities τ_v and K_L completely describe the first order system: K_L describes the steady state performance, and τ_v , the transient. In keeping with mathematical conventions, (II-24) is said to have a first order pole at $-1/\tau_v$.

The linear approximation to a system never reaches a steady state condition but only approaches it asymptotically. Consequently, some figure of merit is needed to describe the speed with which the system approaches the steady-state value. The first order time constant τ is an excellent figure of merit for systems that can be represented by simple transfer functions such as (II-23). For more complex systems a different cri-

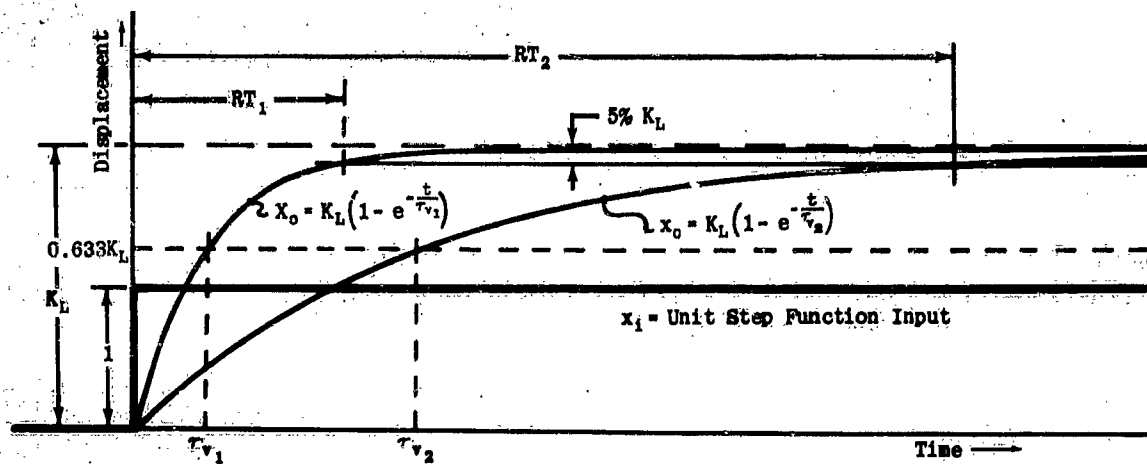


Figure II-61. Response Curve of Hydraulic Amplifier

terion of speed of response is needed. It has been found that a convenient one is the time required for the output to reach and remain within 5% of the final value. This time is referred to by several names such as damping time, settling time, solution time, and response time. By referring to the equation above it can be seen that $x_0(t) \approx 95\% K_L$ when $t = 3\tau$. The response time $RT = 3\tau$. This is shown in figure II-61.

The difference between the time constant (τ) and the response time (RT) is that τ is related strictly to first order poles whereas RT can be related to systems of any degree of complexity. This distinction will become more apparent later.

The preceding discussion has developed the transfer function concept for systems of zero order (G unity) and first order, $G = 1/(\tau s + 1)$. The next simple transfer function of general interest in this sequence is the second order system. The characteristics of this system will be developed in terms of the accelerometer shown in figure II-62.

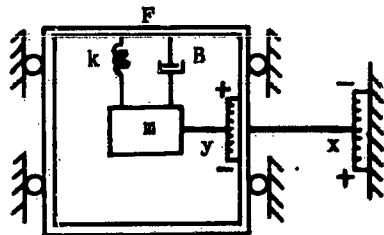


Figure II-62. Accelerometer Assembly

This system consists of a mass (m), a damper (B) and a spring (k). y indicates the motion of the mass relative to the frame and is called the output. x indicates the motion of the frame relative to inertial space and is called the input. The frame (F) is constrained to move in a vertical direction only. Any acceleration of F will cause a displacement of the mass m relative to F , thus giving an indication on the y scale. This is represented more simply in figure II-63 where x_1 indicates motion of mass relative to inertial space.

In order to derive a transfer function relating output

y to input x , the forces on the mass are summed and Newton's law applied. The forces are

Spring Force $F_s = -ky$

Damping Force $F_D = -B(dy/dt)$

Thus:

$F_s + F_D = m(d^2x_1)/(dt^2)$, or $m(d^2x_1)/(dt^2) + B(dy/dt) + ky = 0$.

Replacing x_1 by $y + x$ and rearranging the terms, the equation becomes

$$(II-27) \quad m \frac{d^2y}{dt^2} + B \frac{dy}{dt} + ky = m \frac{d^2x}{dt^2}$$

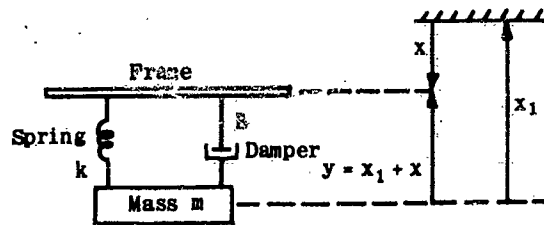


Figure II-63. Equivalent Diagram of Accelerometer

Equation (II-27) is the differential equation of the spring-mass-damper accelerometer assembly. By dividing through by the mass m , the equation is rewritten

$$(II-27a) \quad \frac{d^2y}{dt^2} + \frac{B}{m} \frac{dy}{dt} + \frac{k}{m} y = \frac{d^2x}{dt^2}$$

The coefficient k/m represents the square of the angular undamped natural frequency of the system, ω_n^2 . Likewise B/m can be replaced by the product $2\zeta\omega_n$, where ζ is the damping ratio. Then (II-27a) becomes

$$(II-28) \quad \frac{d^2y}{dt^2} + 2\zeta\omega_n \frac{dy}{dt} + \omega_n^2 y = \frac{d^2x}{dt^2} = a_x$$

where a_x is the acceleration of the frame.

The Laplace transform of (II-28) yields

$$(s^2 + 2\zeta\omega_n s + \omega_n^2) y(s) = a_x(s)$$

Forming the ratio of output to input, (displacement y to

* See any elementary text on dynamics or servomechanisms, e.g., Lauer, Lesnick, and Matson (Ref. 3).

acceleration a_n gives the transfer function of the accelerometer for an acceleration input.

(II-29) denominator in dimensional form denominator in non-dimensional form

$$\frac{Y(s)}{a_n(s)} = \frac{1}{(s)^2 + 2\zeta\omega_n s + \omega_n^2} = \frac{K_a}{\frac{s^2}{\omega_n^2} + \frac{2\zeta}{\omega_n}s + 1}$$

where the gain $K_a = 1/\omega_n^2$. Factoring the dimensional form, (II-29) can be written as:

$$\frac{Y(s)}{a_n(s)} = \frac{1}{(s + \zeta\omega_n + j\omega_n\sqrt{1-\zeta^2})(s + \zeta\omega_n - j\omega_n\sqrt{1-\zeta^2})}$$

Consequently, (II-29) has a pair of complex conjugate poles at $-\zeta\omega_n \pm j\omega_n\sqrt{1-\zeta^2}$.

To determine the nature of the transient for this system, a step acceleration function is applied (by multiplying (II-29) by $1/s$ and the inverse transform obtained. Then

$$y(t) = K_a \left[1 - \frac{1}{\sqrt{1-\zeta^2}} e^{-\zeta\omega_n t} \sin \left(\omega_n \sqrt{1-\zeta^2} t + \tan^{-1} \frac{\sqrt{1-\zeta^2}}{\zeta} \right) \right]$$

This response is plotted in figures II-64 and II-65 for several values of ζ . Each of the curves in the figure is typical of a range of values of ζ . The systems described by such curves are classified as follows:

- $\zeta > 1$ Overdamped second order system.
- $\zeta = 1$ Critically damped second order system.
- $0 < \zeta < 1$ Underdamped second order system.
- $\zeta = 0$ Zero damped (neutrally stable) second order system.

Second order systems with $0 < \zeta < 1$ are of the greatest importance to the designer, so a special series of plots is contained in figure (A-1). * (Note that output is plotted versus the non-dimensional time parameter t/T_n where T_n is the undamped natural period: $T_n = 2\pi/\omega_n$).

These plots reveal that the systems with values of ζ between 0.6 and 0.7 have the lowest response time.

- * Since many of the illustrations used in this chapter are useful in actual engineering problems, they have been placed in the appendix for more convenient everyday reference.

As a matter of fact, it can be shown that the second order system with $\zeta \approx 0.64$ has the least response time for a fixed ω_n . Consequently, in many problems this is a criterion for design.

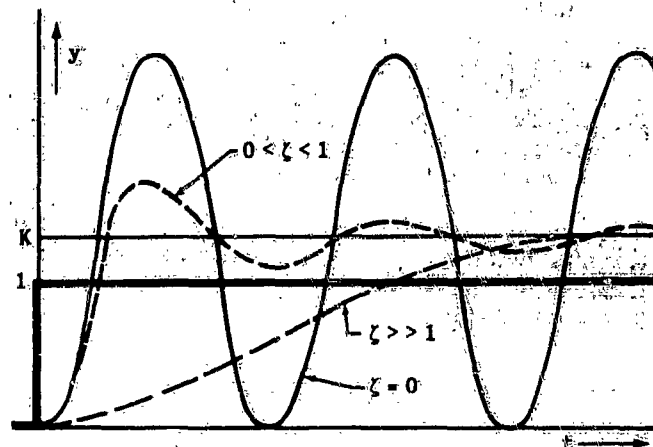


Figure II-64. Three Types of Transient Response to Step Function Input for Three Ranges of the Damping Ratio ζ

The value of ζ in a second order system transfer function must often be determined from experimentally obtained curves. Figures A-1, A-2, A-3, and A-4 are useful for this purpose. * When ζ is close to unity, the overshoots are not clearly enough distinguishable to apply these charts. Consequently, a special plot is used, figure A-5.

In the particular case where $\zeta = 1$, (II-29) takes the form

$$\frac{Y(s)}{a_n(s)} = \frac{1}{(s + \frac{1}{\tau})^2} = \frac{K_a}{(\tau s + 1)^2}$$

where $K_a = \tau^2$ and $\tau = 1/\omega_n$. The inverse transform of this system when a step acceleration function input ($1/s$) is applied is $y(t) = K_a [1 - (1 + t/\tau)e^{-t/\tau}]$. This is plotted in figure II-65.

- * These curves are self explanatory. However, special note must be made that these curves are approximations and must be used with caution.

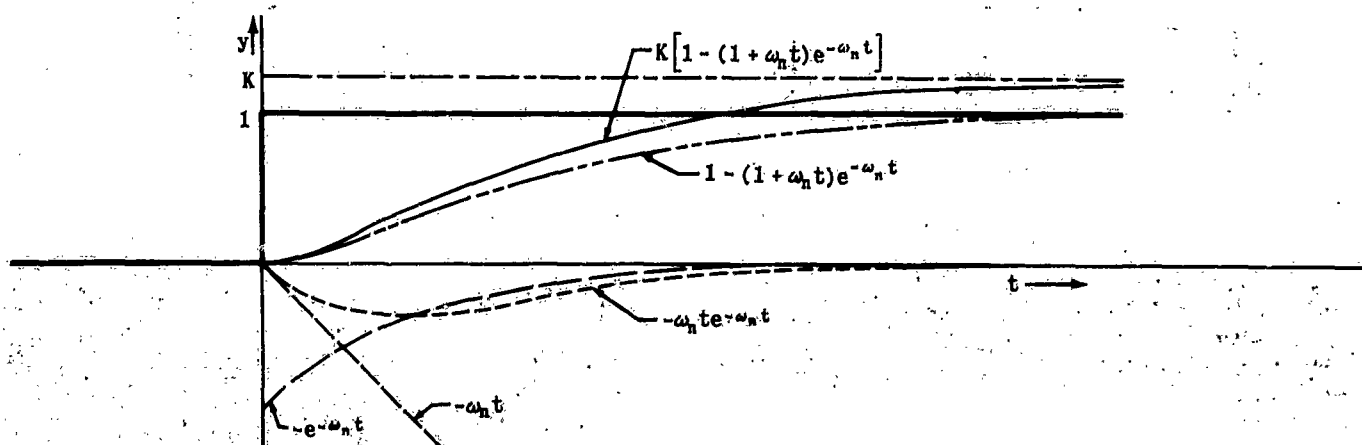


Figure II-65. Second Order System Response with Critical Damping ($\zeta = 1$)

Chapter II

Section 3

When $\zeta > 1$, the characteristic equation (denominator of the transfer function in equation II-29) can be factored into two first order terms:

$$\frac{y(s)}{a_x(s)} = \frac{1}{\left(s + \frac{1}{\tau_1}\right)\left(s + \frac{1}{\tau_2}\right)} = \frac{K_a}{(\tau_1 s + 1)(\tau_2 s + 1)}$$

where $\tau_1 = \frac{1}{\omega_n[\zeta - \sqrt{\zeta^2 - 1}]}$, $\tau_2 = \frac{1}{\omega_n[\zeta + \sqrt{\zeta^2 - 1}]}$, and $K_a = \tau_1 \tau_2$.

The time response of this system to a step acceleration function input is shown in figure II-64. It has the form of the sum of two first order transients.

It is evident from this discussion that second order transfer functions can be written in any of several forms depending on the value of ζ . Defining $\nu = \tau_1/\tau_2$ (see figure A-6):

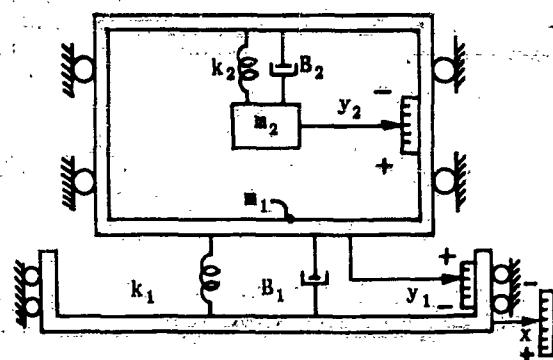
$$(II-30) \quad \frac{y(s)}{a_x(s)} = \frac{1}{s^2 + 2\zeta\omega_n s + \omega_n^2}$$

$$(II-31) \quad \frac{y(s)}{a_x(s)} = \frac{1}{\left(s + \frac{1}{\tau}\right)^2} \quad (\zeta = 1)$$

$$(II-32) \quad \frac{y(s)}{a_x(s)} = \frac{1}{\left(s + \frac{1}{\tau_1}\right)\left(s + \frac{1}{\tau_2}\right)} \quad (\zeta > 1)$$

$$(II-33) \quad \frac{y(s)}{a_x(s)} = \frac{1}{s^2 + \frac{2}{\tau} s + \frac{1}{(\zeta\tau)^2}}$$

$$(II-34) \quad \frac{y(s)}{a_x(s)} = \frac{1}{s^2 + \left(\frac{\nu+1}{\tau_1}\right)s + \frac{\nu}{\tau_1^2}} = \frac{1}{\left(s + \frac{1}{\tau_1}\right)\left(s + \frac{\nu}{\tau_1}\right)} \quad (\zeta > 1)$$



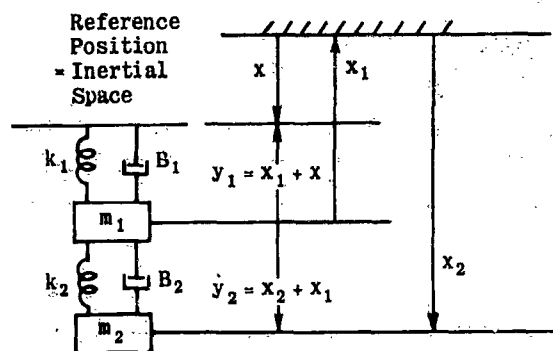
(a)

the various forms of first and second order system transfer functions and their associated time responses. The table includes several forms not discussed above. It is to be noted, for example, that the forms $1/(-\tau s + 1)$ and $1/(s^2 - 2\zeta\omega_n s + \omega_n^2)$ correspond to systems whose time responses increase in value with time. The divergent systems arise from the fact that the exponential factors are positive. All of the mathematical manipulations performed above apply equally well to these diverging systems except that proper accounting of signs must be kept.

Another point of interest is that multiplying the numerator or denominator of the transfer function by s has no effect on the convergence or divergence of the system. Only the initial and final values of the time responses are affected.

The preceding lengthy discussion on the first and second order systems was aimed toward developing a familiarity with their many, essentially equivalent, forms and meanings. The most important fact to be understood is that the gains, poles, and zeros of the transfer functions completely define the behavior of the linear system represented. The charts and tables referred to in the appendix are intended to help the control system designer determine and interrelate these parameters readily. Since nearly all linear system transfer functions can be reduced to products of first and second order factors, even the most complex system can be handled by means of the charts. The following examples will bring out some of these points.

The preceding discussion concerned itself with single



(b)

Figure II-66. Accelerometer Mounted on Spring Damper Arrangement

The relationships between all these parameters are summarized in Table A-1. Another table, A-2, summarizes methods of determining these parameters by measuring certain characteristics of the transient responses.

The third table in the Appendix (Table A-3) summarizes

degree of freedom systems. A two degree of freedom case will now be discussed in order to illustrate some additional features of transfer functions. In this case, two independent transfer functions are required to define the system completely. In figure II-66 the accelerometer is mounted on another spring-damper arrangement. Summing up the forces on each mass, the following equations are obtained:

$$(II-35) \quad m_2 \frac{d^2 y_2}{dt^2} + B_2 \frac{dy_2}{dt} + k_2 y_2 - m_2 \frac{d^2 y_1}{dt^2} = -m_2 \frac{d^2 x}{dt^2}$$

$$B_2 \frac{dy_2}{dt} + k_2 y_2 + m_1 \frac{d^2 y_1}{dt^2} + B_1 \frac{dy_1}{dt} + k_1 y_1 = m_1 \frac{d^2 x}{dt^2}$$

By rearranging the terms and Laplace transforming, the equations are rewritten

$$(II-36) \quad (m_2 s^2 + B_2 s + k_2) y_2(s) - m_2 s^2 y_1(s) = -m_2 s^2 x(s)$$

$$(B_2 s + k_2) y_2(s) + (m_1 s^2 + B_1 s + k_1) y_1(s) = m_1 s^2 x(s)$$

From these equations, the following matrix equation is formed

$$\begin{bmatrix} (m_2 s^2 + B_2 s + k_2) & -m_2 s^2 \\ (B_2 s + k_2) & (m_1 s^2 + B_1 s + k_1) \end{bmatrix} \begin{bmatrix} y_2 \\ y_1 \end{bmatrix} = \begin{bmatrix} -m_2 s^2 \\ m_1 s^2 \end{bmatrix} x$$

from which

$$\frac{y_2(s)}{x(s)} = \frac{\begin{vmatrix} -m_2 s^2 & -m_2 s^2 \\ m_1 s^2 & m_1 s^2 + B_1 s + k_1 \end{vmatrix}}{\begin{vmatrix} m_2 s^2 + B_2 s + k_2 & -m_2 s^2 \\ B_2 s + k_2 & m_1 s^2 + B_1 s + k_1 \end{vmatrix}}$$

and

$$\frac{y_1(s)}{x(s)} = \frac{\begin{vmatrix} m_2 s^2 + B_2 s + k_2 & -m_2 s^2 \\ B_2 s + k_2 & m_1 s^2 \end{vmatrix}}{\begin{vmatrix} m_2 s^2 + B_2 s + k_2 & -m_2 s^2 \\ B_2 s + k_2 & m_1 s^2 + B_1 s + k_1 \end{vmatrix}}$$

Expanding the determinants gives the transfer functions relating outputs y_1 and y_2 to input x .

$$(II-37) \quad \frac{y_2(s)}{x(s)} = \frac{s^2 [-m_2 (m_1 s^2 + B_1 s + k_1) + m_1 m_2 s^2]}{(m_1 s^2 + B_1 s + k_1) (m_2 s^2 + B_2 s + k_2) + m_2 s^2 (B_2 s + k_2)}$$

$$(II-38) \quad \frac{y_1(s)}{x(s)} = \frac{s^2 [m_1 m_2 s^2 + (m_1 + m_2) B_2 s + (m_1 + m_2) k_2]}{(m_1 s^2 + B_1 s + k_1) (m_2 s^2 + B_2 s + k_2) + m_2 s^2 (B_2 s + k_2)}$$

Collecting terms

$$(II-39) \quad \frac{y_2(s)}{x(s)} = \frac{-\frac{m_2}{k_2} s^2 \left(\frac{B_1}{k_1} s + 1 \right)}{As^4 + Bs^3 + Cs^2 + Ds + E}$$

$$(II-40) \quad \frac{y_1(s)}{x(s)} = \frac{\frac{m_1 + m_2}{k_1} s^2 \left[\frac{m_1 m_2}{(m_1 + m_2) k_2} s^2 + \frac{B_2}{k_2} s + 1 \right]}{As^4 + Bs^3 + Cs^2 + Ds + E}$$

where

$$A = \frac{m_1 m_2}{k_1 k_2}$$

$$B = \frac{m_1 B_2}{k_1 k_2} + \frac{m_2 B_1}{k_2 k_1} + \frac{m_2 B_2}{k_2 k_1}$$

$$C = \frac{B_1 B_2}{k_1 k_2} + \frac{m_1}{k_1} + \frac{m_2}{k_2} + \frac{m_2}{k_1}$$

$$D = \frac{B_1}{k_1} + \frac{B_2}{k_2}$$

$$E = 1$$

For this particular system, let

$$m_1 = 100 \text{ slugs}, \quad m_2 = 0.5 \text{ slug}$$

$$k_1 = 5 \text{ lb/ft}, \quad k_2 = 3 \text{ lb/ft}$$

$$B_1 = 1 \text{ lb/ft/sec}, \quad B_2 = 2 \text{ lb/ft/sec}$$

In order to determine the transfer function of the output (y_2) of the accelerometer versus the input (x), these numbers are substituted into (II-39), yielding

$$(II-41) \quad \frac{y_2(s)}{x(s)} = \frac{-0.167 s^2 (2s + 1)}{3.33s^4 + 13.43s^3 + 20.40s^2 + .87s + 1}$$

The fourth order characteristic equation is factorable into quadratics as follows:

$$(II-42) \quad \frac{y_2(s)}{x(s)} = \frac{-0.167 s^2 (2s + 1)}{(.166s^2 + .668s + 1)(20s^2 + .20s + 1)}$$

$$\frac{y_2(s)}{x(s)} = \frac{-0.167 s^2 (2s + 1)}{\left[\left(\frac{s}{\omega_{n1}} \right)^2 + 2\zeta_1 \left(\frac{s}{\omega_{n1}} \right) + 1 \right] \left[\left(\frac{s}{\omega_{n2}} \right)^2 + 2\zeta_2 \left(\frac{s}{\omega_{n2}} \right) + 1 \right]}$$

The first quadratic factor defines a highly damped short period mode ($\zeta_1 = 0.82, \omega_{n1} = 2.46$). The second mode is a poorly damped long period oscillation ($\zeta_2 = 0.023, \omega_{n2} = 0.224$).

For a step acceleration function input ($x(s) = 1/s$), (II-42) becomes

$$(II-43) \quad y_2(s) = \frac{-0.167s(2s + 1)}{(.166s^2 + .668s + 1)(20s^2 + .20s + 1)}$$

The time function $y_2(t)$ is given by the inverse transform of $y_2(s)$. The inverse transform of (II-43) is obtained by a method of partial fraction expansion. By rewriting (II-43) in the dimensional form, it may be shown that

$$(II-44) \quad y_2(s) = \frac{-0.01s(s + 5)}{(s^2 + 4.03s + 6.03)(s^2 + .01s + .05)}$$

$$= \frac{K_1}{(s + \zeta_1 \omega_{n1} - j\omega_{n1} \sqrt{1 - \zeta_1^2})} + \frac{K_2}{(s + \zeta_1 \omega_{n1} + j\omega_{n1} \sqrt{1 - \zeta_1^2})}$$

$$+ \frac{K_3}{(s + \zeta_2 \omega_{n2} - j\omega_{n2} \sqrt{1 - \zeta_2^2})} + \frac{K_4}{(s + \zeta_2 \omega_{n2} + j\omega_{n2} \sqrt{1 - \zeta_2^2})}$$

$$y_2(t) = K_1 e^{(-\zeta_1 \omega_{n1} + j\omega_{n1} \sqrt{1 - \zeta_1^2})t} + K_2 e^{(-\zeta_1 \omega_{n1} - j\omega_{n1} \sqrt{1 - \zeta_1^2})t}$$

$$+ K_3 e^{(-\zeta_2 \omega_{n2} + j\omega_{n2} \sqrt{1 - \zeta_2^2})t} + K_4 e^{(-\zeta_2 \omega_{n2} - j\omega_{n2} \sqrt{1 - \zeta_2^2})t}$$

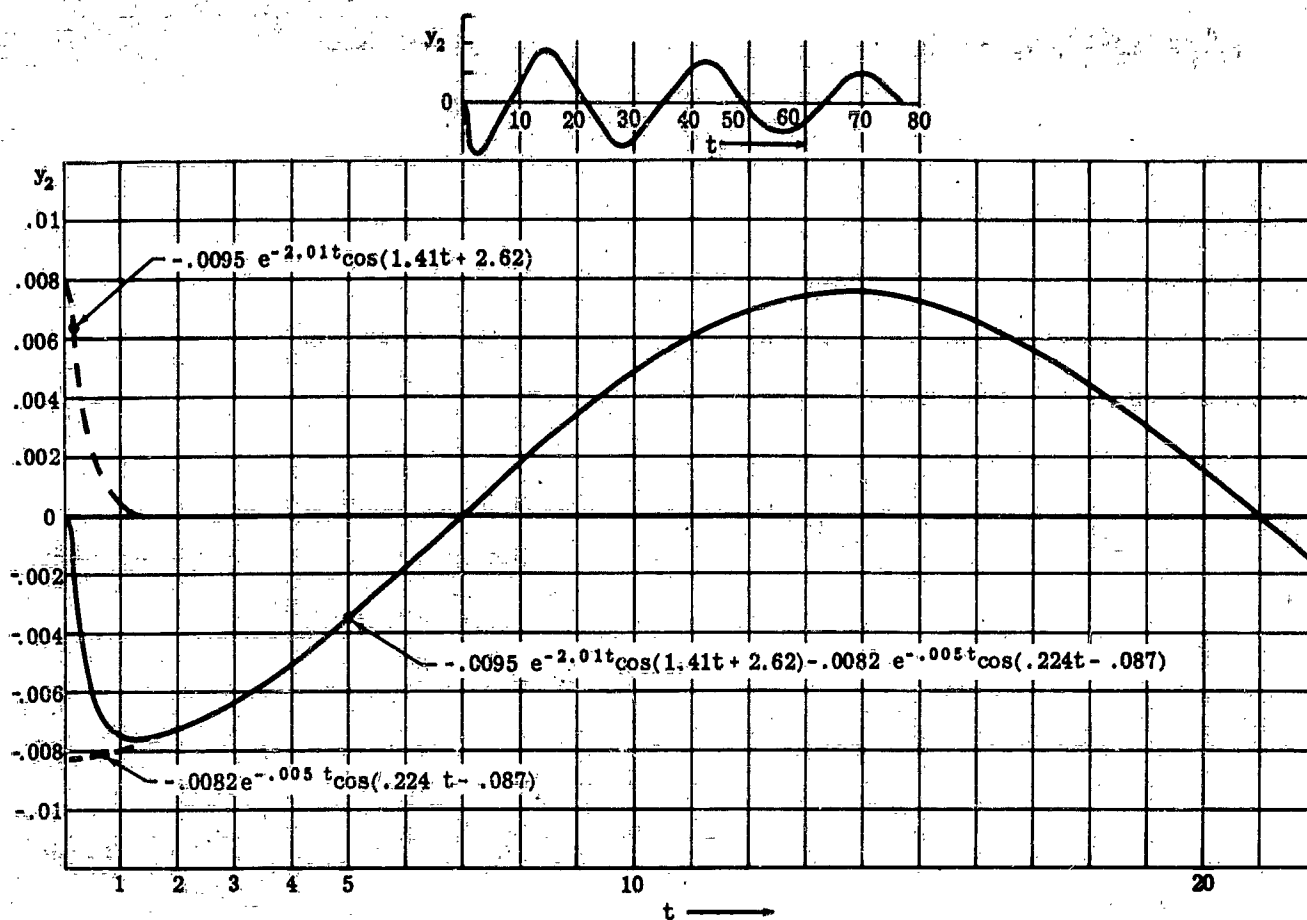


Figure II-67. Transient Response of 2 Degree of Freedom System

where $\zeta_1\omega_{n1} \pm j\omega_{n1}\sqrt{1-\zeta_1^2}$ and $\zeta_2\omega_{n2} \pm j\omega_{n2}\sqrt{1-\zeta_2^2}$ are the roots of the first and second quadratics respectively. The coefficients are*

$$y_2(t) = -0.0095 e^{-2.01t} \cos(1.41t + 2.62) - 0.0082 e^{-0.005t} \cos(.224t - .087)$$

Equation (II-45) is plotted in figure II-67. The first

$$\begin{aligned} K_1 &= \frac{-0.01s(s+5)}{(s + \zeta_1\omega_{n1} + j\omega_{n1}\sqrt{1-\zeta_1^2})(s^2 + 2\zeta_2\omega_{n2}s + \omega_{n2}^2)} \bigg|_{s = -\zeta_1\omega_{n1} + j\omega_{n1}\sqrt{1-\zeta_1^2}} = -0.00475 e^{j2.62} \\ K_2 &= \frac{-0.01s(s+5)}{(s + \zeta_1\omega_{n1} - j\omega_{n1}\sqrt{1-\zeta_1^2})(s^2 + 2\zeta_2\omega_{n2}s + \omega_{n2}^2)} \bigg|_{s = -\zeta_1\omega_{n1} - j\omega_{n1}\sqrt{1-\zeta_1^2}} = -0.00475 e^{-j2.62} \\ K_3 &= \frac{-0.01s(s+5)}{(s + \zeta_2\omega_{n2} + j\omega_{n2}\sqrt{1-\zeta_2^2})(s^2 + 2\zeta_1\omega_{n1}s + \omega_{n1}^2)} \bigg|_{s = -\zeta_2\omega_{n2} + j\omega_{n2}\sqrt{1-\zeta_2^2}} = -0.0041 e^{-j.087} \\ K_4 &= \frac{-0.01s(s+5)}{(s + \zeta_2\omega_{n2} - j\omega_{n2}\sqrt{1-\zeta_2^2})(s^2 + 2\zeta_1\omega_{n1}s + \omega_{n1}^2)} \bigg|_{s = -\zeta_2\omega_{n2} - j\omega_{n2}\sqrt{1-\zeta_2^2}} = -0.0041 e^{j.087} \end{aligned}$$

The final solution is given by

$$(II-45) \quad y_2(t) = C_1 e^{-\zeta_1\omega_{n1}t} \cos(\omega_{n1}\sqrt{1-\zeta_1^2}t + \phi_1) + C_2 e^{-\zeta_2\omega_{n2}t} \cos(\omega_{n2}\sqrt{1-\zeta_2^2}t + \phi_2)$$

* See Gardner and Barnes Page 154 (Ref. 5), for complete discussion of this method.

transient mode is seen to die out very quickly because of its high damping ratio and natural frequency, while the second transient decays very slowly because of the low damping ratio. Since it is a simple matter to obtain the response time of the envelope of an exponentially decaying oscillation from the single parameter $\zeta\omega_n$, and a relatively difficult one to obtain the response time of the actual damped wave, the following approximate re-

relationships are often used:

$$\begin{aligned} 0 < \zeta < 1 & RT \approx \frac{3}{\zeta \omega_n} \\ \zeta \approx 1 & RT \approx \frac{5}{\omega_n} \\ \zeta > 1 & RT \approx \frac{3}{\omega_n} \end{aligned}$$

Thus the response time for the first mode may be approximated by $RT_1 \approx 3/(\zeta_1 \omega_{n1}) = 3/2.01 = 1.49$ seconds. That is, the first transient mode decays to within approximately 5% of its final value in 1.49 seconds. However, the poorly damped mode requires 600 seconds to die down to approximately 5% of its final steady state value $RT_2 \approx 3/(\zeta_2 \omega_{n2}) = 3/0.005 = 600$ seconds.

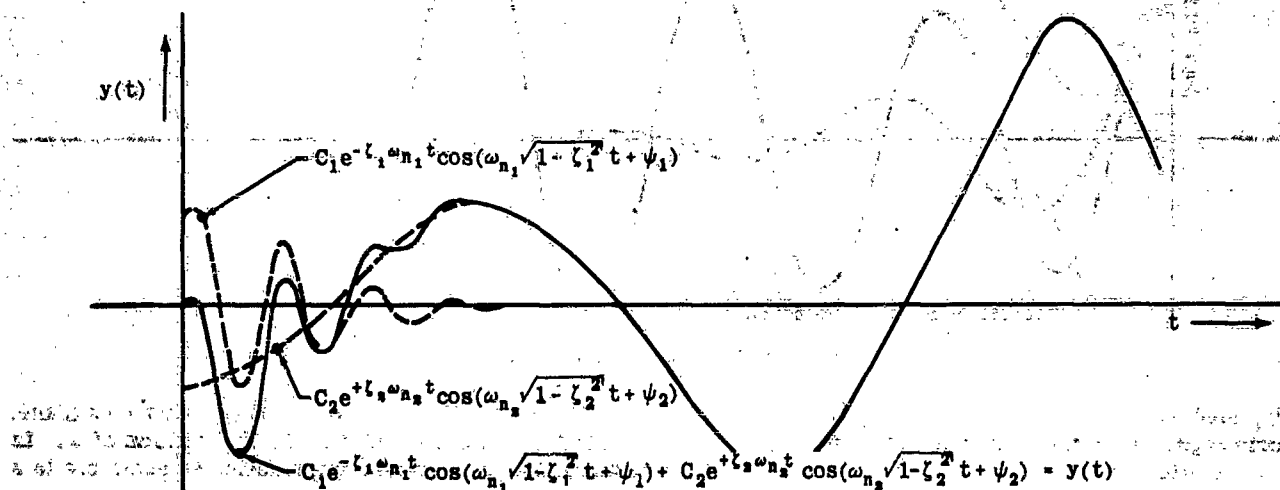


Figure II-68. Response of 2 Degree of Freedom System.

Although it is not physically possible for the system of figure II-66 to be described by such a transform, it is of interest for illustrative purposes to examine an equation of the form of (II-46).

$$(II-46) \quad Y(s) = \frac{K s(s + a_0)}{(s^2 + 2\zeta_1 \omega_{n1} s + \omega_{n1}^2)(s^2 - 2\zeta_2 \omega_{n2} s + \omega_{n2}^2)}$$

The second quadratic factor has a root with a positive real part, leading to a divergent exponential $e^{+\zeta_2 \omega_{n2} t}$. The time domain solution is

$$(II-47) \quad y(t) = C_1 e^{-\zeta_1 \omega_{n1} t} \cos(\omega_{n1} \sqrt{1 - \zeta_1^2} t + \psi_1) + C_2 e^{+\zeta_2 \omega_{n2} t} \cos(\omega_{n2} \sqrt{1 - \zeta_2^2} t + \psi_2)$$

The two modes are given by items 15 and 11, respectively, of Table A-3. A representative plot of (II-47) is shown in figure II-68 with $\zeta_1 > \zeta_2$, $\omega_{n1} > \omega_{n2}$ and $C_1 \approx C_2$.

The diverging oscillation is an example of an unstable mode. The word stability has been avoided up to this point because it means different things to different people. In this volume: If a temporary change in the input to a system causes a temporary change in the output, the system is said to be stable. It is important to understand that this definition says nothing about the detailed behavior of the system. Thus it may approach steady state conditions in a jerky or unsteady manner, but as long as it reaches a steady state it is stable.

In linear systems, the manner in which a system subsides or diverges is indicated by the conventions illustrated by figure II-69. The diagram shows that:

1. The statically stable system tends to return to equilibrium after being displaced (like an ordinary pendulum).
2. The statically unstable system tends to diverge away from equilibrium (like an inverted pendulum).
3. The dynamically stable system returns to equilibrium if it is statically stable.
4. The dynamically unstable system tends to return to equilibrium but it overshoots, reverses direction, and overshoots an even larger amount and thus continues to oscillate at an ever increasing amplitude.

5. The statically unstable, dynamically unstable system has a tendency to diverge away from equilibrium while oscillating with ever increasing amplitude.

All of these characteristics are revealed by the transfer function. Referring to table A-3, it will be observed that when $\zeta < 0$, the system has a dynamically unstable mode while $\zeta > 0$ indicates dynamic stability. A first order term with negative τ indicates static instability and a positive τ , static stability.

In all the transfer function examples, it should be noted that the transfer function is made up of ratios of products of first and second order terms. In fact, the transfer function of any system described by a higher order differential equation can be expressed as a ratio of products of first and second order terms s , $(\tau s + 1)$, and $[s^2/\omega_n^2 \pm (2\zeta/\omega_n)s + 1]$.

An examination of the transfer functions derived so far reveals that they are consistently of the form $Y(s) = Kf(s)$. That is, there is a constant multiplied by a non-dimensional function of s . The transfer functions $Y(s)$ will usually occur as elements in a closed-loop system. The transfer function of units in the forward path of a system is designated by the notation $K_f G$,

where K_g is a positive constant, and G is a non-dimensional ratio of polynomials $N(s)/D(s)$ such that $\lim_{s \rightarrow 0} N(s) = \lim_{s \rightarrow 0} G(s) = +1$. That is, $G(s)$ is of the form

$$\frac{(\pm\tau_1 s + 1) \left(\frac{s^2}{\omega_{n1}^2} + \frac{2\zeta_1}{\omega_{n1}} s + 1 \right) \dots}{(\pm\tau_2 s + 1) (\pm\tau_3 s + 1) \left(\frac{s^2}{\omega_{n2}^2} + \frac{2\zeta_2}{\omega_{n2}} s + 1 \right) \dots}$$

methods to be used make extensive use of graphical constructions involving transfer functions. The characteristics of the graphical representations of the transfer function are discussed in this section.

Basically, the transfer function may be represented in two ways. The first is a plot of the singular values

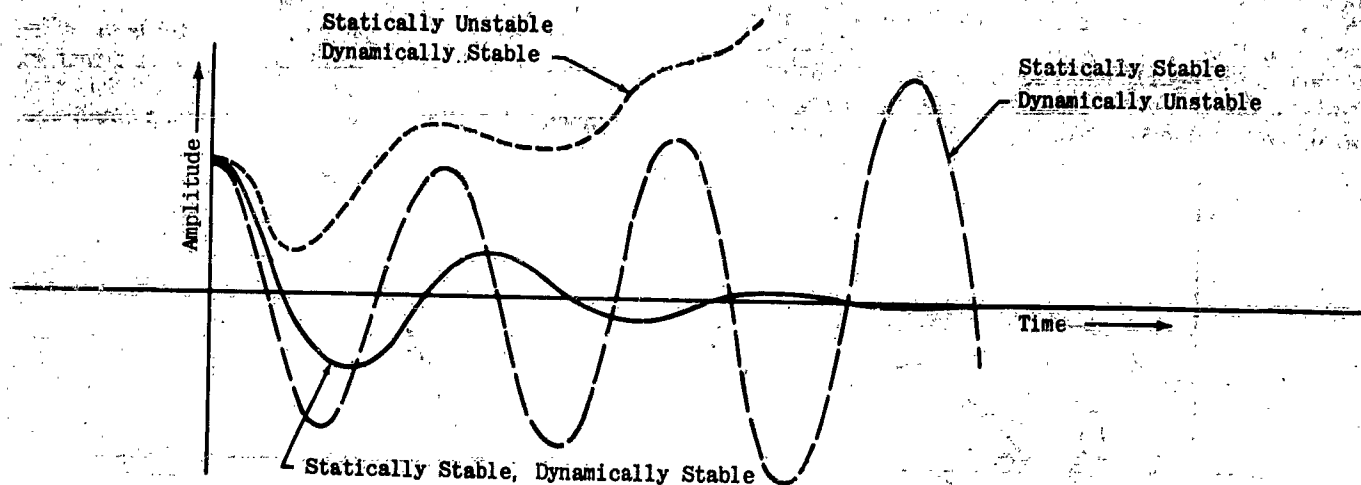


Figure II-69. Stability Curves

Similarly, elements in the feedback path are designated by the form $\pm K_h H$. Thus the closed-loop transfer function can be written

$$\frac{C}{R} = \frac{K_g G}{1 + K_g K_h G H} = \frac{1}{K_h H} \frac{K_g K_h G H}{1 + K_g K_h G H} = \frac{1}{K_h H} \frac{Y(s)}{1 + Y(s)}$$

where $Y(s) = K_g K_h G H$.

It is a comparatively simple task to obtain G and H in factored form as shown, consequently $Y(s)$ is available in factored form and is easily interpreted in terms of performance in the time domain. However, when the operation $1 + Y(s)$ is performed the factors are lost. This is illustrated by the case of a simple servo in which the simple form $Y(s) = K_g / [s(\tau_n s + 1)]$ becomes $Y(s) / [1 + Y(s)] = 1 / [(\tau_n / K_g) s^2 + (1/K_g) s + 1]$. In the first expression, K_g is a simple gain factor, while in the second expression it affects the roots. Consequently, if it were necessary to work exclusively with the second form, the denominator would have to be factored every time the gain K_g were adjusted. In control system design, equations of very high order are common and this process would be very time consuming. Thus, all the methods of control system analysis are directed toward determining performance by working with $Y(s)$ instead of $Y(s) / [1 + Y(s)]$.

(e) GRAPHICAL FORMS OF THE TRANSFER FUNCTION

The preceding section has shown that the system equation or transfer function defines the system performance in terms of its zeros and poles. The problem of analysis then resolves itself into one of determining the zeros and poles of the closed-loop equation. The

of s (the poles and zeros) as discrete points on a plane. The second is a plot of $Y(s)$ for all values of s . In either instance, similar information is available in a simple graphical form.

The plot of discrete, singular points is considered first. The singular points of a transfer function occur for values of s equal to the poles and zeros. The plane on which these values are plotted is referred to as the s -plane.

Since the poles and zeros can be complex, pure real, or pure imaginary numbers, the s -plane plot for an illustrative case might appear as in figure II-70. This figure defines all the characteristics of the s -plane plots.

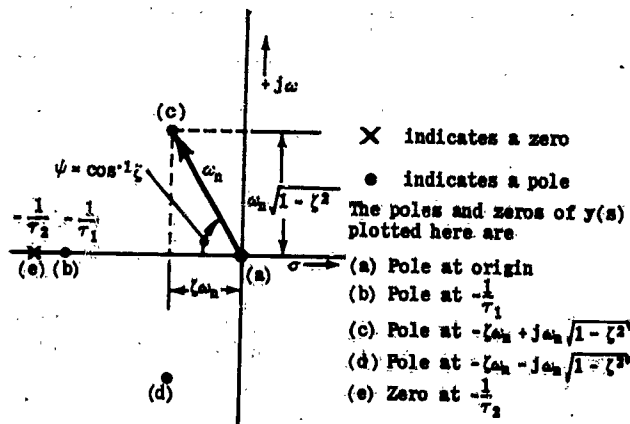


Figure II-70. Poles and Zeros of $Y(s) = \frac{K}{\tau_2 s + 1}$ on the s -Plane

As shown in section II-3d, the poles and zeros determine the type of time response the system will have to an input. Consequently, a table can be constructed relating the s-plane plot to the time response. Table II-2 is such a compilation for a second order system.

Table II-2 indicates that the type of response desired can be controlled by specifying the location of the poles of the closed-loop transfer function. For instance, suppose that for a unit step function input, it is desired that the height of the first overshoot of the output be less than 1.15, i.e. $h/H = .15$. From figure A-2, any value of the damping ratio ζ greater than or equal to .5 will satisfy this condition.

From figure II-71, $\psi = \cos^{-1}\zeta = \cos^{-1}.5 = 60^\circ$. Therefore, excluding all poles of the transfer function from the region to the right of the 60° line as shown in figure II-71 ensures that $\zeta \geq .5$.

If the problem is to make the transient subside to a negligibly small value in a certain time interval, the quantity $\zeta\omega_n = 1/3 RT$ must be controlled. This is done by excluding all poles from the shaded region of figure II-72.

Another important feature is revealed in figure II-70. In particular, a root in the right half of the s-plane

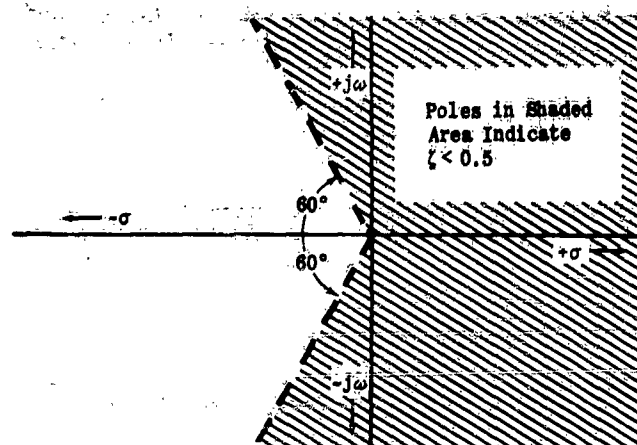


Figure II-71. Poles Excluded from Shaded Region for $\zeta > 0.5$

corresponds to a root with a positive real part. Previous discussions showed that such a root led to a divergent response. This is also indicated in Table II-2. Obviously, then, a requirement for a stable closed-loop system is that it have no poles in the right half plane.

The s-plane is also very useful in determining equation

Damping Ratio	Location of Poles of System Equation in s-Plane	Time Response of Modes
$\zeta < 0$		
$\zeta = 0$		
$0 < \zeta < 1$		
$\zeta = 1$		
$\zeta > 1$		

Table II-2. Effect of Damping Ratio (ζ) on Poles and Transient Response when System

Equation has the Form
$$t(s) = \frac{X_1(s)K}{[(s^2/\omega_n^2) + (2\zeta/\omega_n)s + 1]} \quad \left| \quad X_1(s) = \text{Input} = 1 \right.$$

Chapter II
Section 3

coefficients such as those given by K_1 , K_2 , K_3 , and K_4 in (II-44). Consider a simple closed-loop transfer function given by

$$\frac{C}{R} = \frac{K \kappa}{\left(s + \frac{1}{\tau}\right)(s^2 + 2\zeta\omega_n s + \omega_n^2)} ; \quad \kappa = \frac{\omega_n^2}{\tau}$$

which for a step function input becomes

$$C = \frac{K \kappa}{s \left(s + \frac{1}{\tau}\right)(s^2 + 2\zeta\omega_n s + \omega_n^2)}$$

or

$$C = \frac{K \kappa}{s \left(s + \frac{1}{\tau}\right) (s + \zeta\omega_n - j\omega_n \sqrt{1-\zeta^2}) (s + \zeta\omega_n + j\omega_n \sqrt{1-\zeta^2})}$$

The time solution is given by *

$$C(t) = K_1 + K_2 e^{-\frac{t}{\tau}} + K_3 e^{(-\zeta\omega_n + j\omega_n \sqrt{1-\zeta^2})t} + K_4 e^{(-\zeta\omega_n - j\omega_n \sqrt{1-\zeta^2})t}$$

where

(II-48)

$$K_1 = \frac{K \kappa}{\left(s + \frac{1}{\tau}\right) (s + \zeta\omega_n - j\omega_n \sqrt{1-\zeta^2}) (s + \zeta\omega_n + j\omega_n \sqrt{1-\zeta^2})} \Big|_{s=s_0=0}$$

$$K_2 = \frac{K \kappa}{s \left(s + \zeta\omega_n - j\omega_n \sqrt{1-\zeta^2}\right) (s + \zeta\omega_n + j\omega_n \sqrt{1-\zeta^2})} \Big|_{s=s_1=-\frac{1}{\tau}}$$

$$K_3 = \frac{K \kappa}{s \left(s + \frac{1}{\tau}\right) (s + \zeta\omega_n + j\omega_n \sqrt{1-\zeta^2})} \Big|_{s=s_2=-\zeta\omega_n + j\omega_n \sqrt{1-\zeta^2}}$$

$$K_4 = \frac{K \kappa}{s \left(s + \frac{1}{\tau}\right) (s + \zeta\omega_n - j\omega_n \sqrt{1-\zeta^2})} \Big|_{s=s_3=-\zeta\omega_n - j\omega_n \sqrt{1-\zeta^2}}$$

By making the indicated substitutions for s , (II-48) becomes:

$$(II-49) \quad K_1 = \frac{K \kappa}{\left(\frac{1}{\tau}\right) (\zeta\omega_n - j\omega_n \sqrt{1-\zeta^2}) (\zeta\omega_n + j\omega_n \sqrt{1-\zeta^2})}$$

$$K_2 = \frac{K \kappa}{\left(-\frac{1}{\tau}\right) \left(-\frac{1}{\tau} + \zeta\omega_n - j\omega_n \sqrt{1-\zeta^2}\right) \left(-\frac{1}{\tau} + \zeta\omega_n + j\omega_n \sqrt{1-\zeta^2}\right)}$$

$$K_3 = \frac{K \kappa}{\left(-\zeta\omega_n + j\omega_n \sqrt{1-\zeta^2}\right) \left(-\zeta\omega_n + j\omega_n \sqrt{1-\zeta^2} + \frac{1}{\tau}\right) \left(-\zeta\omega_n + j\omega_n \sqrt{1-\zeta^2} + \zeta\omega_n + j\omega_n \sqrt{1-\zeta^2}\right)}$$

$$K_4 = \frac{K \kappa}{\left(-\zeta\omega_n - j\omega_n \sqrt{1-\zeta^2}\right) \left(-\zeta\omega_n - j\omega_n \sqrt{1-\zeta^2} + \frac{1}{\tau}\right) \left(-\zeta\omega_n - j\omega_n \sqrt{1-\zeta^2} + \zeta\omega_n - j\omega_n \sqrt{1-\zeta^2}\right)}$$

In sub-section (d) the K 's were obtained by substituting the numerical values for each of the parameters (ζ , ω_n , etc.) in (II-49). K_1 is evaluated in this way very easily; K_2 , K_3 , and K_4 somewhat more complicated and graphical computation is advantageous in these cases. This can be done by utilizing the s -plane plot as shown in the following explanation.

The poles ($s=0$, $s_1=-1/\tau$, $s_2=-\zeta\omega_n + j\omega_n \sqrt{1-\zeta^2}$; and $s_3=-\zeta\omega_n - j\omega_n \sqrt{1-\zeta^2}$) of the transfer function are plotted on the s -plane in figure II-73. Also shown are the vectors representing these poles. A vector is drawn from s_2 to s_1 representing the factor

* See Gardner and Barnes, Page 154 (Ref. 5).

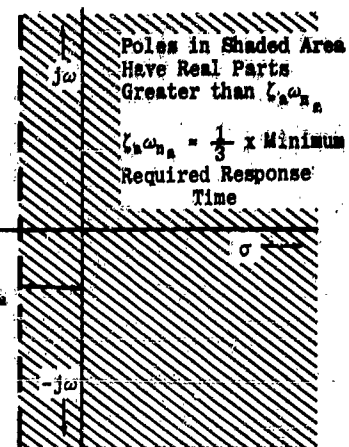


Figure II-72. Area of Exclusion to Obtain Specified Minimum Damping ($\zeta\omega_n$)

$\bar{s}_1 - \bar{s}_2 = (-1/\tau + \zeta\omega_n - j\omega_n \sqrt{1-\zeta^2})$ in equation (II-49). In a similar manner, the vector $(-1/\tau + \zeta\omega_n + j\omega_n \sqrt{1-\zeta^2}) = \bar{s}_1 - \bar{s}_3$ is drawn in figure II-74. This figure shows that $-1/\tau + \zeta\omega_n - j\omega_n \sqrt{1-\zeta^2}$ and $-1/\tau + \zeta\omega_n + j\omega_n \sqrt{1-\zeta^2}$ are complex conjugates. Therefore, the product of these two vectors is a real number:

$$K_2 = \frac{K \kappa}{\left(-\frac{1}{\tau}\right) \left(-\frac{1}{\tau} + \zeta\omega_n - j\omega_n \sqrt{1-\zeta^2}\right)^2}$$

K_2 is obtained by measuring the lengths of two vectors \bar{s}_1 and $\bar{s}_1 - \bar{s}_2$ and carrying out the indicated multiplication $K_2 = (-K\kappa)/(|\bar{s}_1| \cdot |\bar{s}_1 - \bar{s}_2|)^2$

The product is negative and the vector \bar{K}_2 points in the negative direction along the real axis.

Each of the factors of K_3 is shown in figure II-75. K_3 may be expressed as $K_3 = (K\kappa)/(|\bar{K}_3| e^{j\phi_K})$ where $|\bar{K}_3| = |\bar{s}_2| |\bar{s}_2 - \bar{s}_1| |\bar{s}_2 - \bar{s}_3|$ and $\phi_K = \phi_3 + \phi_4 + \phi_5$.

Similarly, \bar{K}_4 is plotted in figure II-76. Note particularly that \bar{K}_4 is the complex conjugate of \bar{K}_3 . So $|\bar{K}_3| = |\bar{K}_4|$ and $\phi_{K_3} = -\phi_{K_4}$.

The complete time solution $C(t)$ is now

$$C(t) = K \left[1 + \frac{\kappa}{\left(-\frac{1}{\tau}\right) \left(-\frac{1}{\tau} + \zeta\omega_n - j\omega_n \sqrt{1-\zeta^2}\right)^2} e^{-\frac{t}{\tau}} + K_3 |e^{(-\zeta\omega_n + j\omega_n \sqrt{1-\zeta^2})t - j\phi_{K_3}}| + K_4 |e^{(-\zeta\omega_n - j\omega_n \sqrt{1-\zeta^2})t - j\phi_{K_4}}| \right]$$

where $|\bar{K}_3| = (K\kappa)/|\bar{K}_3|$ and $|\bar{K}_4| = (K\kappa)/|\bar{K}_4|$. Since \bar{K}_3 and \bar{K}_4 are complex conjugates

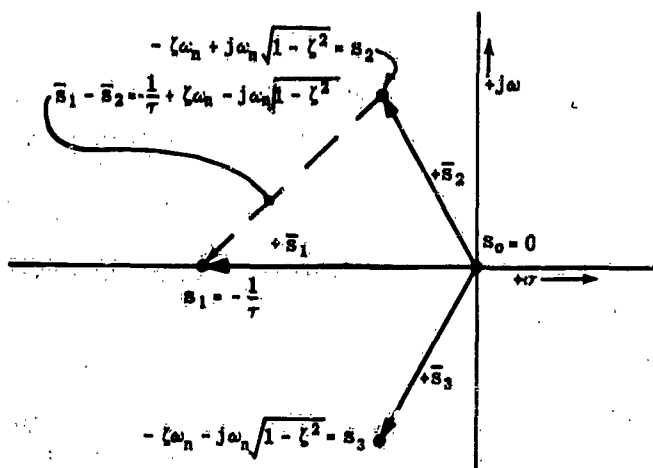


Figure II-73. Vector Plot of $C = \frac{K K}{s(s+\frac{1}{T})(s^2+2(\omega_n s + \omega_n^2))}$

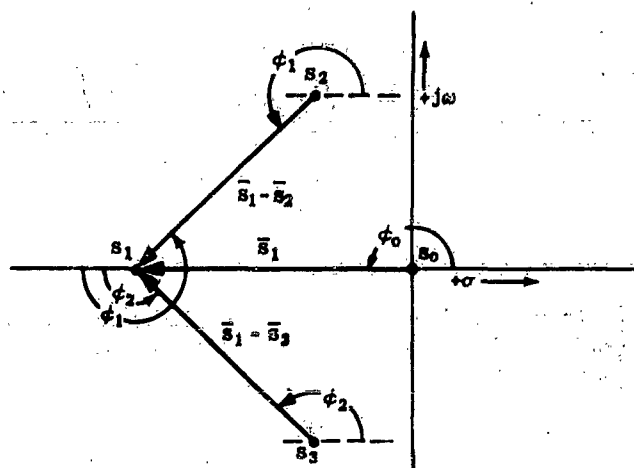


Figure II-74. Graphical Solution for K_2

$$\begin{aligned} & |K_3| e^{-j\phi_{K_3} + (-\zeta\omega_n + j\omega_n\sqrt{1-\zeta^2})t} + |K_4| e^{-j\phi_{K_4} + (-\zeta\omega_n - j\omega_n\sqrt{1-\zeta^2})t} = \\ & |K_3| [e^{-j\phi_{K_3} + (-\zeta\omega_n + j\omega_n\sqrt{1-\zeta^2})t} + e^{j\phi_{K_3} + (-\zeta\omega_n - j\omega_n\sqrt{1-\zeta^2})t}] = \\ & |K_3| e^{-\zeta\omega_n t} [e^{j(\omega_n\sqrt{1-\zeta^2} - \phi_{K_3}')} + e^{-j(\omega_n\sqrt{1-\zeta^2} - \phi_{K_3}')}] \end{aligned}$$

Referring to figure II-77 it is seen that the term

$$[e^{-j(\phi_{K_2} - \omega_n \sqrt{1-\zeta^2}t)} + e^{j(\phi_{K_2} - \omega_n \sqrt{1-\zeta^2}t)}] = 2 \cos(\omega_n \sqrt{1-\zeta^2}t - \phi_{K_2})$$

is the sum of two unit vectors. Consequently, only figure II-75 need have been constructed to obtain the amplitude terms of the time response.

The time response is:

$$C(t) = K + K_2 e^{-\frac{t}{\tau}} + \frac{2K\kappa}{|\kappa_2|} e^{-\omega_n t} \cos(\omega_n \sqrt{1-\zeta^2} t - \phi_{\kappa_2})$$

Note that this value checks with the inverse transform given on page 343 of Gardner and Barnes (Ref. 5). It

is interesting to observe that the vectors K_3' and K_4' fall in the left half of the s-plane. Because of this, it might be assumed that the amplitude of the last term of the time response should be preceded by a negative sign. However, this would be incorrect since the sign is taken care of in the term $\cos(\omega_n \sqrt{1-\zeta^2} t - \alpha_3)$.

This example had no zeros in the transfer function. To illustrate the procedure of graphically determining the transient response amplitude terms for a case in which the transfer function has zeros, the coefficients of (II-44) shall be determined. The graphical solutions for K_1 , K_2 , K_3 and K_4 of (II-44) are given in figures II-78 and II-79. Only K_1 and K_3 are obtained since K_2 is the complex conjugate of K_1 , and K_4 is the complex conjugate of K_3 . The special point to observe in this example is that zeros are treated in the same manner as poles with due caution being taken to keep the algebraic signs correct.

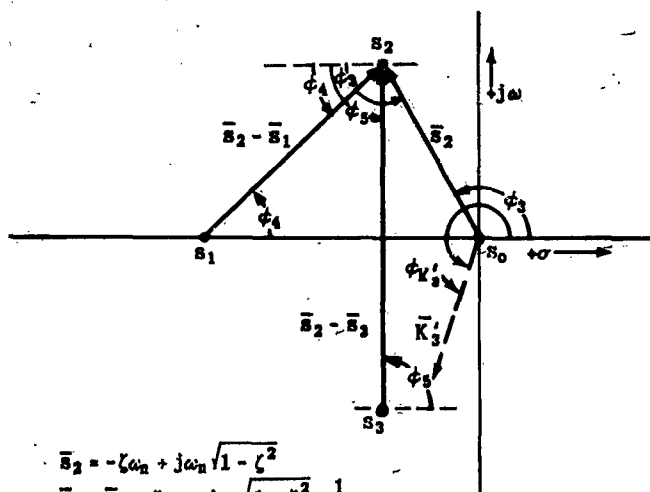


Figure II-75. Graphical Solution for K_2

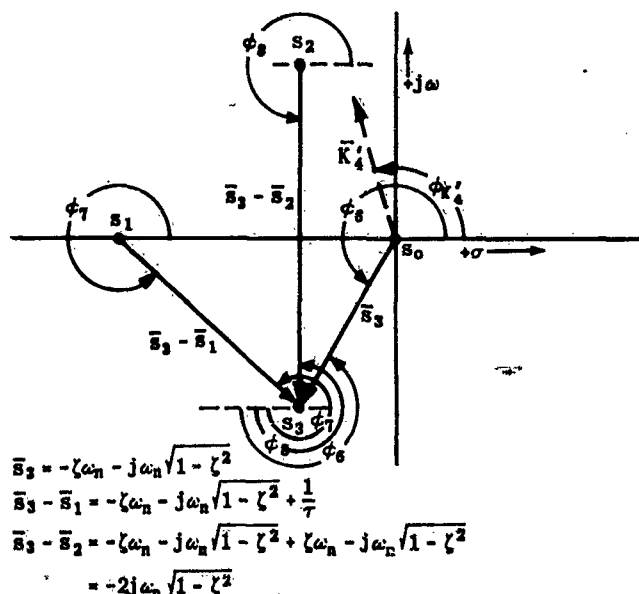


Figure II-76. Graphical Solution for K_4

Chapter II
Section 3

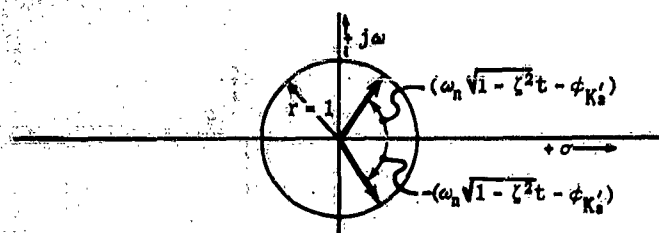


Figure II-77. Conjugate Unit Vectors

Following this procedure, the time solution for (II-44) is obtained by substituting these K 's into the proper equation. Since

$$y_2(t) = K_1 e^{-2.01t + j1.41t} + K_2 e^{-2.01t - j1.41t} + K_3 e^{-.005t + j.224t} + K_4 e^{-.005t - j.224t}$$

substituting the values for the K 's from the graphs, $y_2(t)$ becomes:

$$y_2(t) =$$

$$(-.00475e^{j2.62} e^{-2.01t + j1.41t} - .00475e^{-j2.62} e^{-2.01t - j1.41t}) + (-.0041e^{-j.087} e^{-.005t + j.224t} - .0041e^{j.087} e^{-.005t - j.224t})$$

This may be rewritten as:

$$y_2(t) = -.00475e^{-2.01t} (e^{j(2.62+1.41t)} + e^{-j(2.62+1.41t)}) - .0041e^{-.005t} (e^{j(-.087+.224t)} + e^{-j(-.087+.224t)})$$

or

$$y_2(t) = -.0095e^{-2.01t} \cos(1.41t + 2.62) - .0082e^{-.005t} \cos(.224t - .087)$$

which agrees with (II-45).

The second basic way of representing the transfer function is to plot $Y(s)$ as a continuous function of s . Since s is a complex variable, such a plot for all values of s would require four dimensions. However, very useful representations can be obtained by allowing s to be a pure imaginary ($j\omega$), thereby restricting the plot to two dimensions.

$Y(j\omega)$ is in general a complex number. Its value for any fixed ω can be expressed as a magnitude (called the amplitude ratio) and a phase angle. Two convenient methods of graphical representation will be described in this subsection. In the first, the amplitude ratio is plotted at a given phase angle on a polar

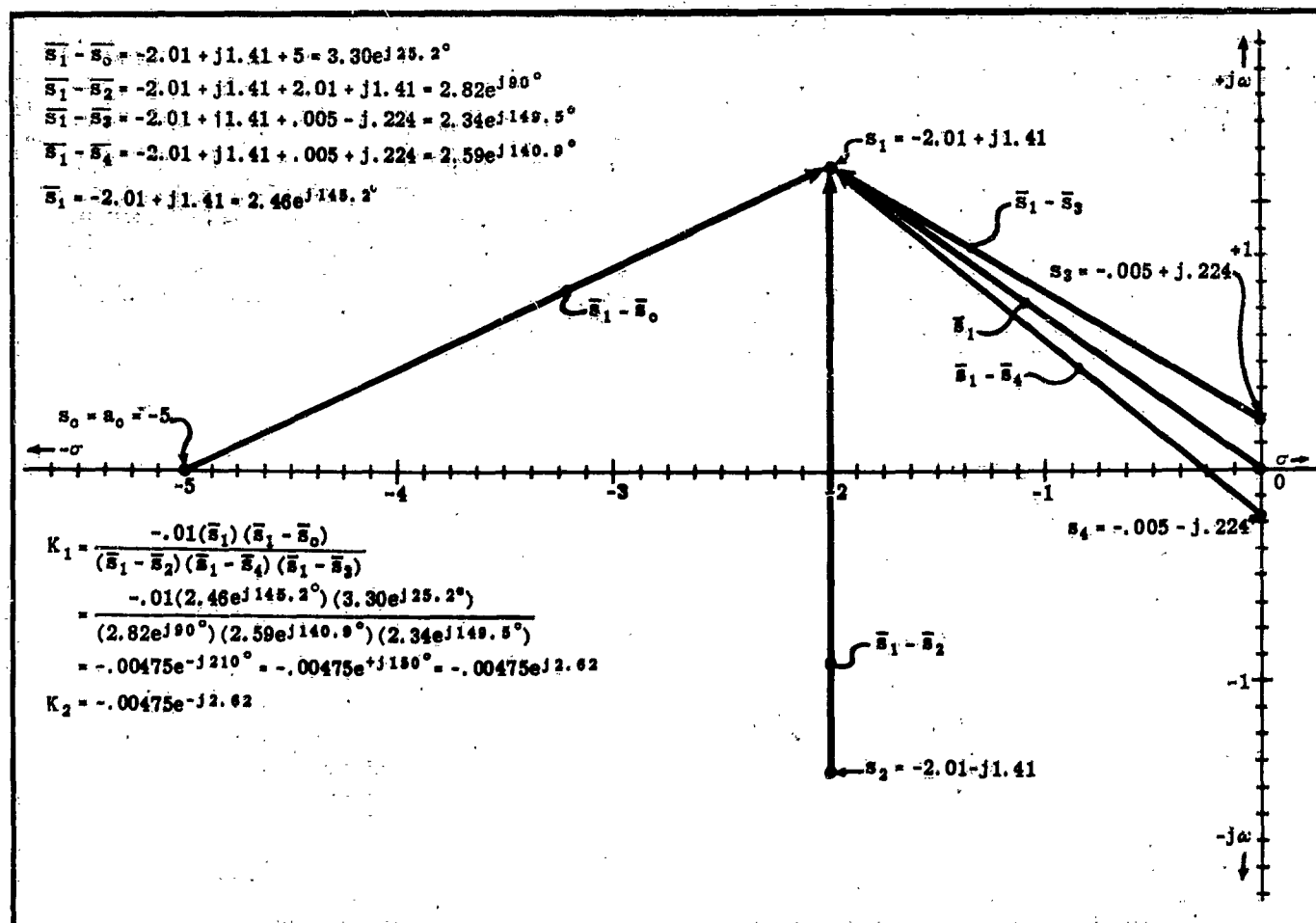
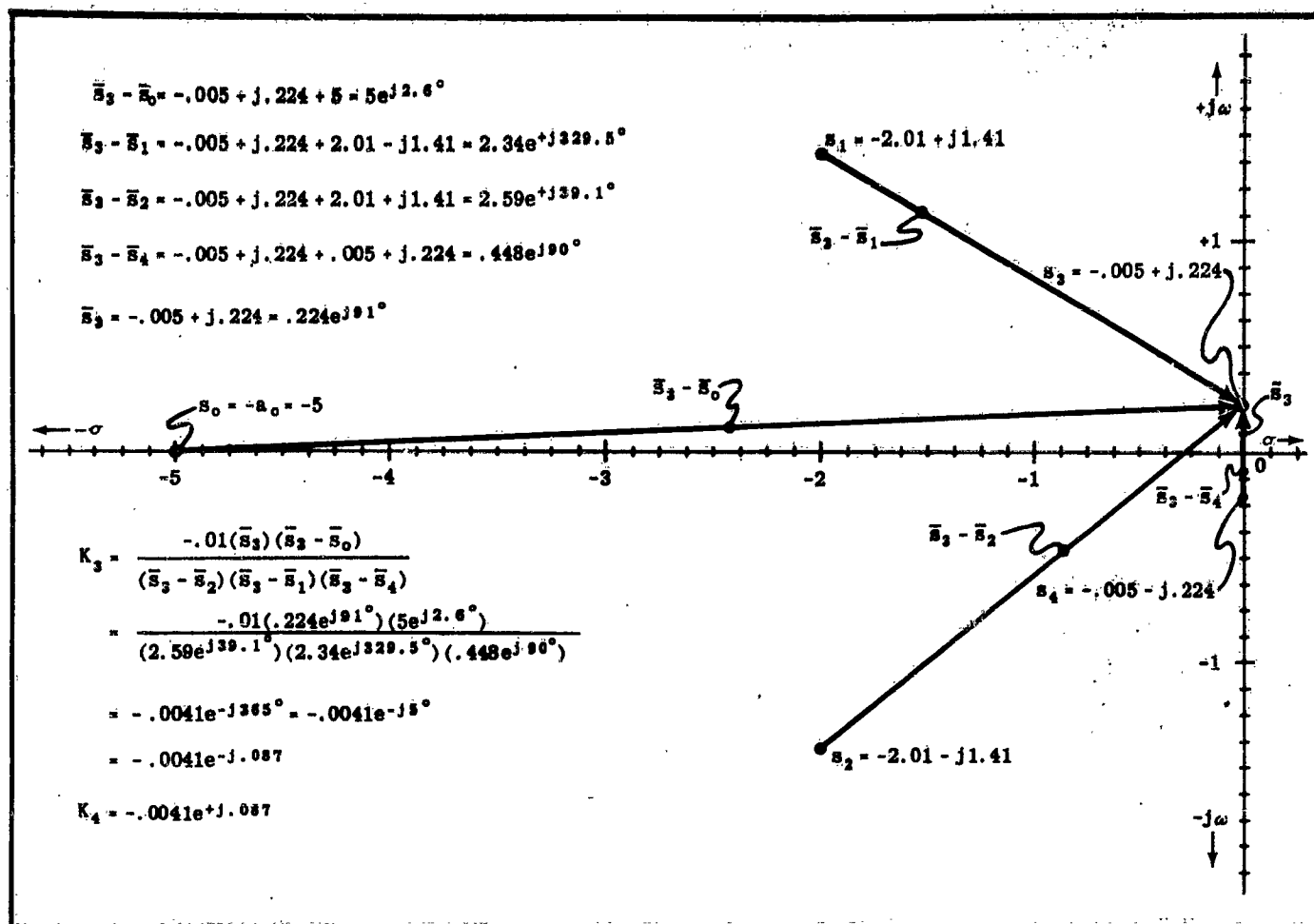


Figure II-78. Graphical Solution For K_1 of (II-44)


 Figure II-79. Graphical Solution for K_3 and K_4 of (II-44)

plot with ω as a parameter. In the second, a convenient function of the amplitude ratio is plotted versus frequency on semi-log paper with the phase angle being plotted in a similar way.

For the purposes of this subsection, the selection of $s = j\omega$ is essentially a matter of convenience. However, there is a definite correlation between $Y(j\omega)$ and part of the time response of a system excited by a sinusoidal input. This correlation will be considered in detail.

In the development of the graphical representations a procedure similar to that used previously will be employed. The explanations will be made for actual examples, starting with very simple cases.

The first graphical representation discussed will be the polar parametric plot commonly referred to as the Nyquist diagram.

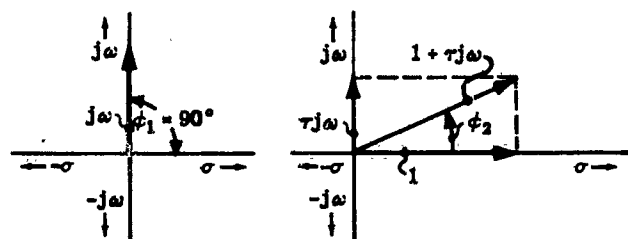
Consider the system which has the transfer function:

$$(II-50) \quad K_g G(s) = \frac{K}{s(\tau s + 1)}$$

If $j\omega$ is substituted for s , where $-\infty < \omega < \infty$,

$$(II-51) \quad K_g G(j\omega) = \frac{K}{j\omega(\tau j\omega + 1)}$$

Figure II-80 shows that both $j\omega$ and $\tau j\omega + 1$ may be represented as vectors.


 Figure II-80. Vector Representations of $j\omega$ and $\tau j\omega + 1$

Writing $K_g G(j\omega)$ in a form to take advantage of this fact, (II-52) is obtained.

$$(II-52) \quad K G(j\omega) = \frac{K}{r_1 e^{j\phi_1} r_2 e^{j\phi_2}}$$

where $r_1 = \omega$, $r = \sqrt{\tau^2 \omega^2 + 1}$, $\phi_1 = 90^\circ$, $\phi_2 = \tan^{-1} \tau \omega$ so that

$$(II-53) \quad K_g G(j\omega) = \frac{K}{\omega \sqrt{\tau^2 \omega^2 + 1}} e^{-j(\phi_1 + \phi_2)}$$

Equation (II-53) shows that $KG(j\omega)$ can be represented as a vector with a magnitude $K/(\omega\sqrt{\tau^2\omega^2 + 1})$ and a phase angle $-(\phi_1 + \phi_2)$. The magnitude is commonly called the amplitude ratio.

Plots of (II-53) for $K = 10$ and $K = 2.5$ with $\tau = .5$ are drawn on a polar chart in figure II-81.

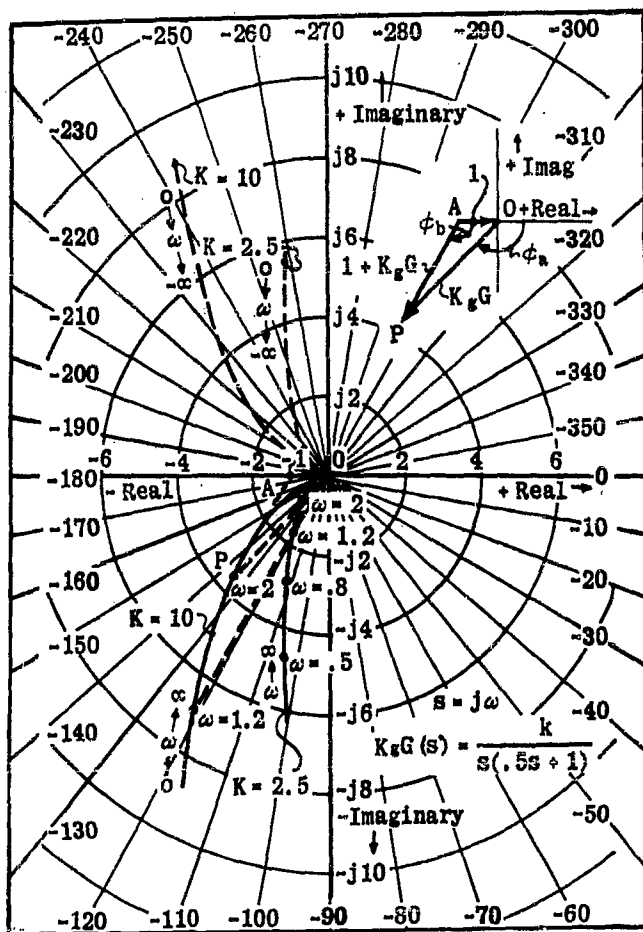


Figure II-81. Nyquist Diagram of Open-Loop Transfer Function

For illustrative purposes the vector for $\omega = 1.2$ and $\omega = 2$ are shown.

Notice that rectangular coordinate axes have been superimposed upon the polar plot with the axis of imaginaries oriented along the $+90^\circ$, $+270^\circ$ line and the axis of reals oriented along the -180° , 0° line. The vector drawn from the point $-1 + j0$ on this coordinate axis to the locus is the vector $1 + K_s G(j\omega)$. Since the closed-loop response is given by

$$(II-54) \quad \frac{C(s)}{R(s)} = \frac{K_s G(s)}{1 + K_s G(s)}$$

$$\frac{C(s)}{R(s)} = \frac{|OP|}{|AP|} \frac{\angle \phi_a}{\angle \phi_b} = \frac{|OP|}{|AP|} \angle \phi_a - \angle \phi_b$$

This relationship makes it possible to obtain the closed-loop magnitude and phase relationships from those of an open-loop Nyquist plot.

Particular attention should be paid to the effect of the constant K in figure II-81. The phase angles are independent of the magnitude of K . Consequently, the shape of the curve is unchanged by changing K . The only effect of K on the polar plot is to change its scale.

Another important feature of the plot is that it is symmetrical about the real axis. The reason for this is that $\angle Y(s)_{s=j\omega}$ changes sign with ω , while the squaring process required to establish $|Y(s)|_{s=j\omega}$ eliminates the effect of a negative " ω ."

Higher order transfer functions are plotted in the same way. The procedure is to set $s = j\omega$ and to calculate the magnitude and phase angle of the transfer function for each value of ω from zero to infinity. Table II-3 is a summary of locus plots for some common transfer functions. In this table only the parts of the loci corresponding to $0 < \omega < +\infty$ are plotted for simplicity.

One part of every Nyquist plot that is of special interest is the "low frequency end." Table II-3 shows that as the frequency, ω , approaches zero certain loci take on infinitely large values and approach an axis as an asymptote (Table II-3d and II-3f). Loci exhibiting this behavior correspond to important classes of control systems. The three most important ones will be discussed in some detail in chapter IV. Figure II-82 shows the characteristic loci for each class and the corresponding form of the transfer function.

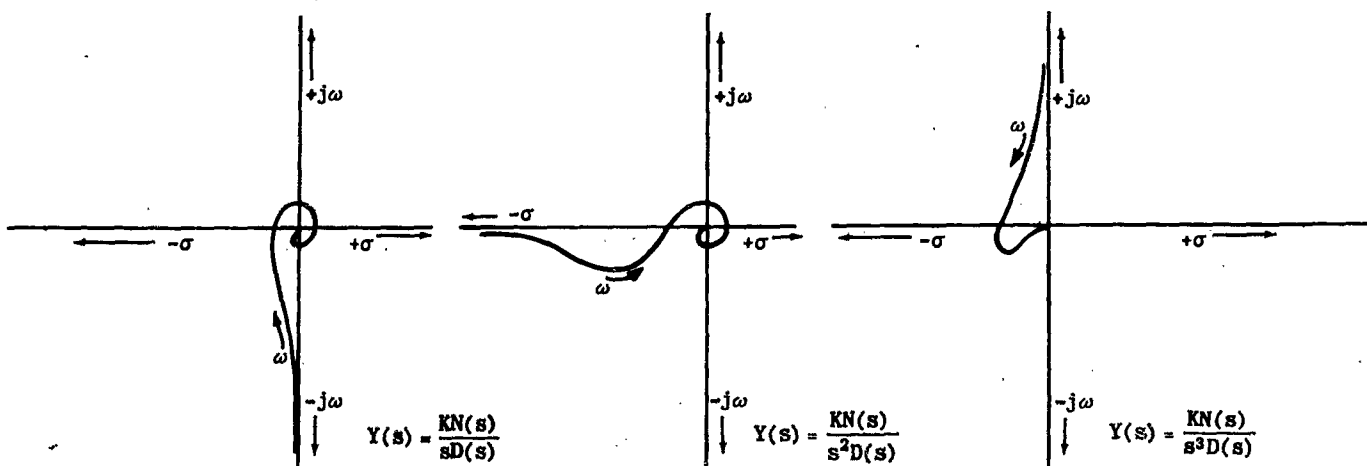


Figure II-82. Special Locus Shapes

In the preceding discussion, the transfer function $|Y(s)|_{s=j\omega}$ has been treated as a convenient abstraction. However, the part of the locus corresponding to $0 < \omega < +\infty$ can be obtained physically by applying a sinusoidal forcing function to the system. The differential equation then becomes of the form

(II-55)

$$a_n \frac{d^n x}{dt^n} + a_{n-1} \frac{d^{n-1} x}{dt^{n-1}} + \dots + a_0 = X_A \sin \omega t, \quad (0 < \omega < +\infty)$$

Transforming,

(II-56) $(a_n s^n + a_{n-1} s^{n-1} + \dots + a_0) X(s) = \frac{X_A \omega}{s^2 + \omega^2}$

$$X(s) = \frac{X_A \omega / a_0}{\left(\frac{a_n}{a_0} s^n + \frac{a_{n-1}}{a_0} s^{n-1} + \dots + 1 \right) (s^2 + \omega^2)}$$

Expressions of this type can be inverse transformed* as follows:

(II-57) $x(t) = \frac{X_A}{a_0} \left\{ i \left[\frac{\frac{a_n}{a_0} s^n + \frac{a_{n-1}}{a_0} s^{n-1} + \dots + 1}{s^2 + \omega^2} \right]_{s=j\omega} \right. \\ \left. + K_1 e^{s_1 t} + K_2 e^{s_2 t} + \dots + K_n e^{s_n t} \right\}$

* See Ref. 5 Page 159

$Y(s)$	Locus on the Complex $Y(s)$ Plane	$Y(s)$	Locus on the Complex $Y(s)$ Plane
(a) $\frac{1}{s}$		(f) $\frac{1}{s \left[\frac{s^2}{\omega_n^2} + \frac{2\zeta}{\omega_n} s + 1 \right]}$	
(b) $\frac{1}{s^2}$		(g) s	
(c) $\frac{1}{\tau s + 1}$		(h) $\frac{\tau_1 s + 1}{\tau_2 s + 1}, \tau_2 > \tau_1$	
(d) $\frac{1}{s(\tau s + 1)}$		(i) $s(\tau s + 1)$	
(e) $\frac{1}{\frac{s^2}{\omega_n^2} + \frac{2\zeta}{\omega_n} s + 1}$		(j) $\frac{s^2}{\omega_n^2} + \frac{2\zeta}{\omega_n} s + 1$	
$\omega \rightarrow$ Indicates Direction of Increasing Frequency			

Table II-3. Some Common Transfer Functions and Their Locus Shapes on the Linear Plot

where I indicates "imaginary part of."

The motion $x(t)$ described by this equation consists of a sinusoidal oscillation of amplitude X_A/a_0 and frequency ω , and a transient determined by the sum of the exponential terms $K_1 e^{s_1 t}$, $K_2 e^{s_2 t}$, etc. The sinusoidal oscillation is described by

$$(II-58) \quad x(t)_{\sin} = I \left[\frac{\frac{X_A}{a_0} e^{st}}{\frac{a_n}{a_0} s^n + \frac{a_{n-1}}{a_0} s^{n-1} + \dots + 1} \right]_{s=j\omega}$$

$$= I \left[\frac{K}{\frac{a_n}{a_0} s^n + \frac{a_{n-1}}{a_0} s^{n-1} + \dots + 1} \right]_{s=j\omega} \sin \omega t$$

where $K = X_A/a_0$ and $0 < \omega \leq +\infty$. The term in the brackets determines the amplitude and phase angle of the sinusoidal response and is recognized as the transfer function $Y(s)|_{s=j\omega}$ of the system. This result applies to linear systems of any order of complexity. It shows that the transfer function can be obtained from the physical system by applying a steady sinusoidal input and separating the steady output oscillation from the transient terms. The stability of a system has no effect on the validity of these results, although, practically speaking it sometimes is difficult to separate the sinusoid from the "transients" in unstable systems without the application of special techniques.

The second graphical representation of $Y(s)$ to be discussed is the logarithmic or Bode plot. The calculations required to construct such a plot point by point are identical to those discussed above. For the simple system

$$(II-59) \quad K_s G(s) = \frac{K}{s(\tau s + 1)}$$

the form (II-53) is derived, and the amplitude ratio $K/(\omega \sqrt{\tau^2 \omega^2 + 1})$ and phase angle $-(\phi_1 + \phi_2)$ are plotted independently versus the angular frequency ω on semi-logarithmic paper. The ω is plotted on logarithmic scale, and $20 \log_{10} [K/(\omega \sqrt{\tau^2 \omega^2 + 1})]$ and $-(\phi_1 + \phi_2)$ are plotted on the linear scale as in figure II-83.

The notation $20 \log_{10}(\quad)$ is referred to as the log-modulus or L_m and performing the indicated operation results in a number in decibel (or db) units (see figure A-20). One advantage of this procedure is that amplitude ratios can be plotted very simply. In the example chosen (see II-51)

$$(II-60) \quad L_m \left| \frac{K}{j\omega(\tau j\omega + 1)} \right| = L_m K - L_m |j\omega| - L_m |\tau j\omega + 1|$$

These three terms are plotted in figure II-83 as curves A, B, and C and are summed as curve D.

The dashed sloping line marked (B) is the $L_m 1/(j\omega)$ curve. Since the plot is done on logarithmic coordinates, it is a straight line as shown. In a plot of this type a change in ω by a factor of 10 is called a "decade;" a change by

a factor of 2 (doubling its value) is referred to as an "octave." Consequently, since the slope of the line is -20 db per decade, the curve is said to "attenuate" at a rate of 20 db/dec. or 6 db/octave. That is, $L_m 1/(j\omega) = 20 \log_{10} 1/\omega = -20 \log_{10} \omega$ and each time the frequency is increased by a factor of 10, the magnitude of $L_m 1/(j\omega)$ is decreased by 20 db. The 20 db/dec. attenuation rate is true of all $s = j\omega$ terms in the denominator of a transfer function. An $s = j\omega$ term in the numerator plots with a positive slope and is said to be amplified at the rate of 20 db/dec. (or 6 db/oct.). The attenuation and amplification rate will be specified in terms of decades in this volume. If the $s = j\omega$ term is raised to the n^{th} power, the slope is $20n$ db/dec., i.e., $L_m |j\omega^n| = 20 \log_{10} \omega^n = 20n \log_{10} \omega$

The $(\tau j\omega + 1)^{-1}$ term of the transfer function is shown by curve (C) in figure II-83. This curve approaches two straight line asymptotes. The relationships for establishing the asymptotes to the true curve are as follows: when $j\omega\tau \ll 1$,

$$(II-61) \quad L_m \left| \frac{1}{j\omega\tau + 1} \right| \sim L_m 1 = 0 \text{ db}$$

When $j\omega\tau \gg 1$,

$$(II-62) \quad L_m \left| \frac{1}{j\omega\tau + 1} \right| \sim L_m \left| \frac{1}{j\omega\tau} \right|$$

Equation (II-61) establishes a horizontal line at zero db and (II-62) establishes a straight line sloping at -20 db per decade. The asymptotes intersect at the frequency where $\omega = 1/\tau$; ($\omega\tau = 1$). At $\tau\omega = 1$, $L_m 1/(j\omega\tau + 1) = L_m 1/\sqrt{2} \approx -3$ db. Therefore, the true curve lies -3 db from the asymptote at this point. At one octave below the "break point," $\tau\omega = .5$. Therefore $L_m 1/(j\omega\tau + 1) = L_m 1/(\sqrt{.5^2 + 1}) = L_m 1/(\sqrt{1.25}) \approx -1$ db. The actual log modulus value is then 1 db below the asymptote. At $\tau\omega = 2$, (one octave above the break point), $L_m 1/(\sqrt{2^2 + 1}) = L_m 1/\sqrt{5} \approx -7$ db, however, at $\tau\omega = 2$, the asymptote is at -6 db. Therefore, the log-modulus curve lies one db below the asymptote at this point.

The complete log-modulus curve, (D), for the transfer function (II-59) is obtained by simply adding the three curves (A), (B), and (C) in accordance with (II-60). The same result is also obtained by adding the asymptotes for each of the factors of the transfer function, line $b o a'$, and then sketching in the true curve $b c a'$. Notice that the break point is still at $\omega = 2$, ($\tau\omega = 1$) and that the 3 db and 1 db departure characteristics still hold. For any changes in the value of K the db scale need only be shifted up or down depending on the nature of the change of K . Evidently this procedure is general for first order terms and may be summarized as follows:

1. Establish break point ($\omega = 1/\tau$) on the zero db line and draw in the asymptotes.
2. At the break point, draw a line sloping downward to the right at 20 db per decade for a denominator term, and upward to the right for a numerator term.
3. At the break point, spot a point 3 db below the zero db line for a denominator term and above the line for a numerator term.
4. At one octave above and at one octave below

the break point, spot points 1 db distant from the asymptotes, below for a denominator term and above for a numerator term.

5. Sketch in the log-modulus curve, using the asymptotes and the three points.

The phase angle curves for the transfer function are also shown in figure II-83; the phase angle varies from -90° to -180° as the frequency increases (solid line). This phase angle curve is the sum of the two curves arising from the $(j\omega)^{-1}$ and $(.5j\omega + 1)^{-1}$ factors. The $(j\omega)^{-1}$ term plots as a constant -90° phase angle

(curve E). The $(.5j\omega + 1)^{-1}$ angle factor approaches an asymptote at 0° and another at -90° ; its midpoint (45°) being determined by $1/\tau$. Figure A-7 is included in the appendix in order to facilitate sketching the curves. It is to be noted that a factor $\tau s + 1$ in the numerator produces a phase curve starting at a 0° asymptote and approaching $+90^\circ$ at high frequencies and of the same shape as the denominator factor.

Many of these same principles can be used in plotting second order terms in logarithmic form. However, since the second order factor is a function of two in-

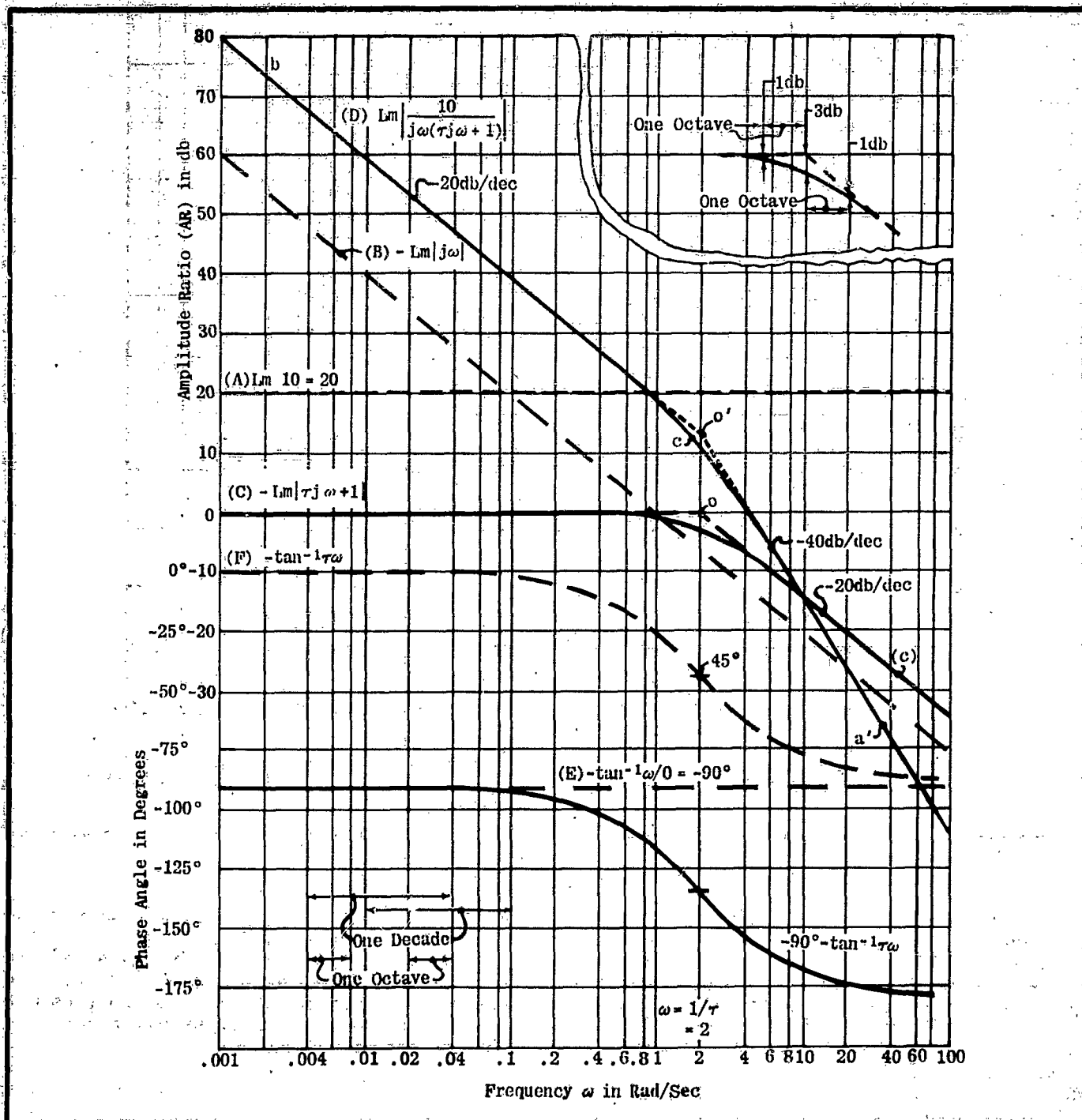


Figure II-83. Bode Diagram of Open Loop Transfer Function $K G(s) = \frac{K}{s(\tau s + 1)}$, $K = 10$, $\tau = .5$

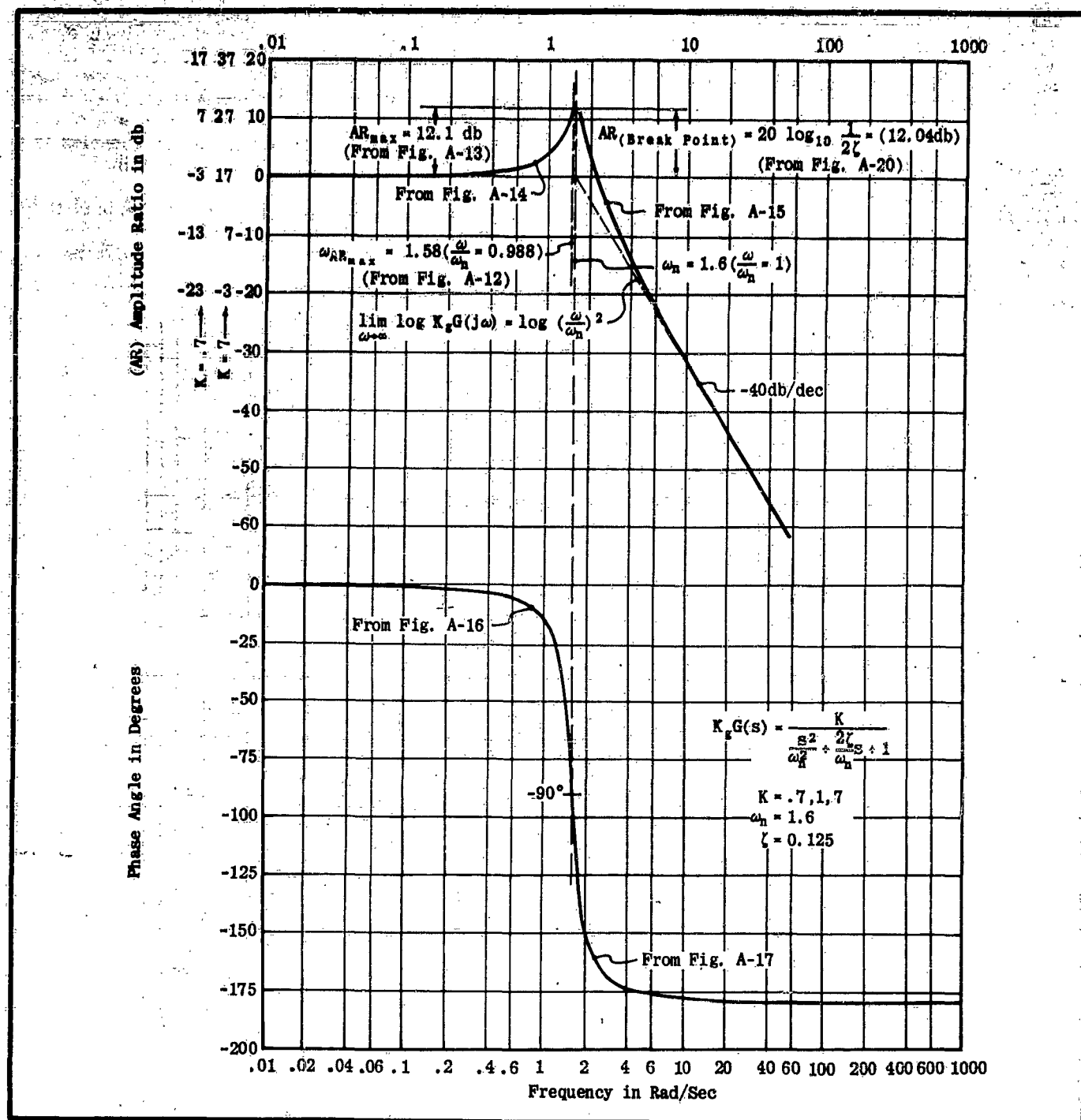


Figure II-84. Bode Diagram of $K_r G(s)$.

dependent parameters ζ and ω_n , plotting is slightly more complicated. Figure II-84 points out some of the special features of these terms. The two key characteristics are the -40 db/dec slope at high frequencies and the phase shift of 180°.

Several charts have been prepared to aid in the construction of Bode plots of second order factors. After the break point and the -40 db/dec asymptote have been established, figures A-12 and A-13 are used to locate the peak when one exists. Figures A-14 through A-17 are charts which give amplitude ratio departure from the asymptotes and phase angle. For most purposes

a satisfactory sketch of the amplitude curve can be completed using only figures A-12 and A-13. However, figures A-14 and A-15 can be used to aid in the construction of more accurate curves. Values of the amplitude departure from the asymptotes for discrete frequency ratios ω/ω_n are read by proceeding up and down the ordinate representing the value of ζ . The same procedure is followed when using figures A-16 and A-17 to construct the phase angle plot.

Gain (K_r) adjustments are made by shifting the db scale as shown in figure II-84. Consequently, the amplitude curves are nearly always plotted for

$K_f = 1$ (0db), and the 0 db line corresponding to the true gain marked in later. If the gain (K_f) is greater than unity, the true 0 db line occurs below the one used for plotting; a K_f less than unity places it above.

It is clear that any number of factors of any order may be added on Bode plots to achieve a complete transfer function plot. Since it is so simple to plot first and second order terms, the transfer functions are always factored accordingly. This avoids tedious computations of amplitudes and phase angles of complicated transfer functions.

Table A-4 is a summary of forms of the transfer functions encountered in system analysis. In the last column showing the Bode plots of the factors, only the asymptotes to the log-modulus curves are shown. For items (7, 8, 9, and 10), the 3 db and 1 db departure relationships apply. For the second order terms (items 11, 12, 13, and 14), the exact shape of the phase angle curve and the log-modulus (amplitude ratio) curves depend on the damping ratio ζ .

Notice particularly in Table A-4 that the asymptote curve breaks upward for all numerator terms (items 1, 2, 3, 7, 9, 11, and 13), and all denominator terms show a downward break of the asymptote. Note also that this fact plus the phase change indicate whether the zeros or poles are in the right or left half of the s-plane. That is, a phase curve tending to go in the same direction as the amplitude ratio curve (items 7, 8, 11, and 12) shows that the zeros or poles are in the left half plane, while those that go in the opposite direction indicate zeros or poles in the positive half of the s-plane.

Bode (Ref. 8) refers to systems that contain no poles or zeros in the right half plane as minimum phase systems. If any poles or zeros exist in the right half plane, the system is non-minimum phase. Following this lead, those factors that represent poles or zeros in the right half s-plane are generally referred to as non-minimum phase terms; all others, including poles and zeros on the imaginary axis, are minimum phase terms.

The following general conclusions can be made concerning the interpretation of Bode diagrams:

1. Asymptote slopes must always be either zero or some integral multiple of ± 20 db/dec.
2. The change in slope of the asymptotic plot at a break point indicates the order of the pole or zero that exists at the break point.
3. A positive change in slope corresponds to a zero.
4. A negative change in slope corresponds to a pole.
5. a. The location of the break point of the asymptotes indicates the location of first order poles and zeros on the s-plane.
b. The location of the break point and the departure from the asymptotes indicates location of second order poles and zeros on the s-plane.
6. When phase and amplitude curves change in the same direction, a minimum phase term is indicated.
7. When phase and amplitude curves change in opposite directions a non-minimum phase term occurs.
8. When the slope of the amplitude as $\omega \rightarrow 0$ is -20 db/dec. the system is of the zero position error type.*
9. When the slope is -40 db/dec. as $\omega \rightarrow 0$ the system is of the zero velocity error type.*
10. When the slope is -60 db/dec. as $\omega \rightarrow 0$ the system is of the zero acceleration error type.*

SECTION 4 - SERVOMECHANISMS

The control system field is extraordinarily broad and most of the previously discussed methods are sufficient to describe any linear problems in this field. However, this book is concerned with the methods of handling only a certain class of control systems.

From the discussion in the preceding pages it is evident that there are two broad classes of control systems: Open loop control systems and feedback control systems. There are aircraft flight control systems in both of the classes. Typical examples of an open loop system are the common cable or push-pull rod surface controls; on the other hand, a hydraulic valve cylinder combination or an autopilot are closed loop systems.

Open loop systems are by their very nature calibrated systems, and their performance is profoundly affected by the condition of the calibration. Aircraft flight control system designers are intimately aware of this and expend a great deal of time attempting to minimize environmental effects on the calibration by means of such devices as cable tension regulators. In any case, the design principles governing those devices are well known. Although one of the end results of these volumes will be to set up criteria governing the performance of these systems, the design process itself in these cases

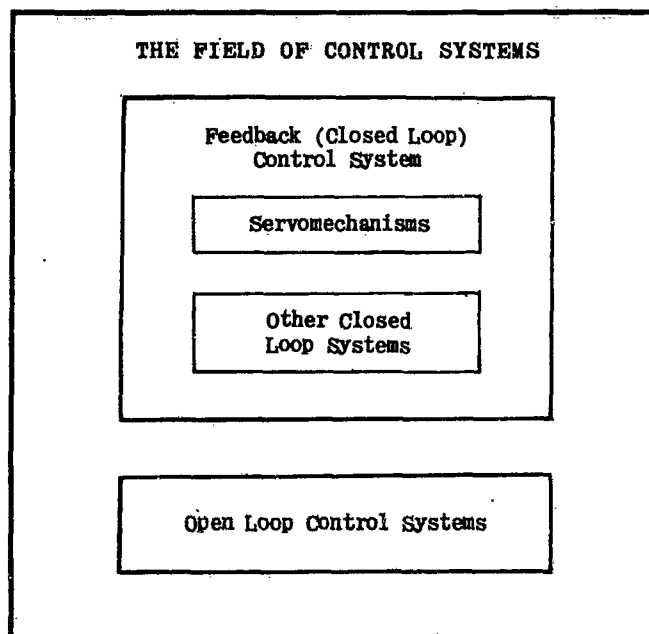


Figure II-85. The Field of Control Systems

* See Chapter IV.

Servomechanisms have the following defining characteristics:

1. They are closed loop systems.
2. A large amount of power is controlled by a relatively low power element.

3. They involve mechanical motion.

This is a very broad definition and includes many systems not often brought to mind by common parlance. Thus such things as autopilots, tracking control systems, pilot-airframe combinations, etc., are referred to as servomechanisms or servo systems.

This classification of control systems is illustrated in figure II-93.

BIBLIOGRAPHY

The following bibliography is included for reference. The list is in no sense complete, but contains the major source material for this chapter. Many of the references, themselves, contain much more complete and detailed bibliographies.

1. 'Cybernetics, or Control and Communication in the Animal and Machine', by N. Wiener; John Wiley and Sons, New York, 1948.
2. 'Human Use of Human Beings; Cybernetics and Society,' by N. Wiener; Houghton Mifflin Co., 1950.
3. 'Servomechanism Fundamentals,' by Lauer, Lesnick, and Matson; McGraw Hill Book Co., New York, 1947.
4. 'Theory of Servomechanisms,' by James, Nichols, and Phillips; McGraw Hill Book Co., New York, 1947.
5. 'Transients in Linear Systems,' by M. F. Gardner and J. L. Barnes; John Wiley and Sons, New York, 1942.
6. 'Block Diagram Network Transformation,' by T. D. Graybeal; Electrical Engineering, Vol. 70, No. 11, November, 1951.
7. 'Regeneration Theory,' by H. Nyquist; Bell System Tech. Jour., Vol. 11, January 1932.
8. 'Network Analysis and Feedback Amplifier Design,' by H. W. Bode; D. Van Nostrand Co., New York, 1945.

* In this volume, the control systems to be considered are feed back (closed-loop) systems, and, particularly, those known as servomechanisms.

CHAPTER III

ANALYSIS

SECTION 1 — INTRODUCTION

System analysis is concerned with an inquiry into the behavior of a given system. Previous chapters have established that the static and dynamic performance of a linear system, including its responses to known inputs, is completely determined by its transfer function. The problem of linear analysis then becomes one of obtaining information about system transfer functions. Since servo analysis is primarily concerned with feedback control systems, the linear analysis problem is further limited to obtaining information regarding closed loop transfer functions from a knowledge of the open loop transfer functions. This chapter will consider the important methods and techniques available for solving this limited problem.

The notion of a feedback control system, such as that represented by the block diagram of figure III-1, has been previously introduced and the algebra of such block diagrams has been considered. With all of this background information understood, the essential problem of linear servo analysis can be described.

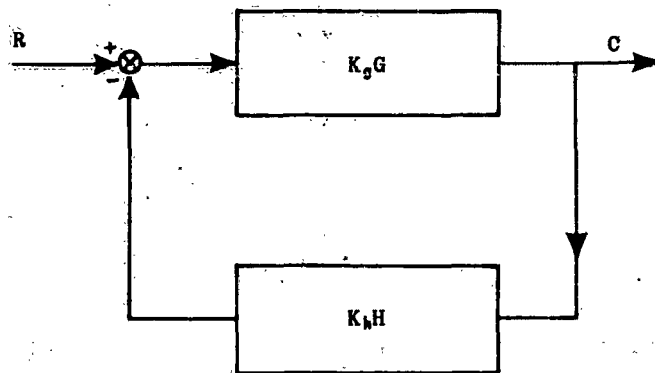


Figure III-1. Illustrative Servomechanism

In the control system represented by figure III-1, the closed loop transfer function, which defines the properties of the system, is given by:

$$(III-1) \quad \frac{C}{R} = \frac{1}{K_H H} \left[\frac{K_G K_H G H}{1 + K_G K_H G H} \right] = \frac{1}{K_H H} \left[\frac{Y}{1 + Y} \right]$$

where $Y = K_G K_H G H$. The analysis problem is essentially solved when the properties of the closed loop transfer functions are known to the analyst. Transfer functions are completely specified by their poles, zeros, and scale factors. Therefore, the analysis problem to be con-

sidered is concerned with gathering information about the values of the poles and zeros of the closed loop transfer function $(1/K_H H)[Y/(1 + Y)]$ from a knowledge of the open loop transfer function, Y . Since $K_H H$ is known, the portion of (III-1) requiring additional study is the bracketed term, $Y/(1 + Y)$.

While the poles and zeros of $Y/(1 + Y)$ are the prime information required, the major effort of analysis need be directed toward finding only the poles, since the zeros of $Y/(1 + Y)$ are the zeros of Y , and hence, are known. To illustrate, if $N(s)$ is the numerator of Y and $D(s)$ the denominator,

$$(III-2) \quad Y = \frac{N(s)}{D(s)}$$

$$(III-3) \quad \frac{Y}{1 + Y} = \frac{N(s)/D(s)}{1 + N(s)/D(s)} = \frac{N(s)}{N(s) + D(s)}$$

The analysis problem can now be stated mathematically as: Given Y , determine the poles of $Y/(1 + Y)$, or alternatively the zeros of $1 + Y$.

Before the specific content of the chapter is outlined, it should be mentioned that a direct analytical method of determining the poles and zeros is to factor the closed loop transfer function. However, for all but the simplest systems, this procedure may be very tedious and time consuming. Therefore, direct factorization is usually impractical and will not be discussed in this chapter. However, methods of approximate factorization are included in an appendix to this volume, and may be used if desired.

The major portion of this chapter consists of three interrelated sections. These sections are organized so that the techniques employed give closed loop zero and pole locations with greater and greater accuracy as one method succeeds another.

The first method presented enables one to obtain only very general information concerning the regions in which the poles and zeros lie. In addition to this information, a certain amount of qualitative data can sometimes be obtained by analogy between the behavior of actual systems and very simple systems by the use of transfer function characteristics discussed in the fifth section of this chapter. A rule known as the Generalized Cauchy-Nyquist criterion is developed and used as the basis of this method. Nyquist diagram and s-plane representations of the transfer function

are utilized in the application of this criterion.

The second method permits a much more exact determination of closed loop pole and zero values. While the previous method requires the use of both Nyquist diagram and s-plane plots, this section utilizes the logarithmic transfer function representation and another graphical aid called the Nichols chart.

The third method is the most exact presented. The values of closed loop poles and zeros are determined to accuracies limited only by the graphical process involved. Only the s-plane representation of the transfer function is required, from which the loci of closed loop

poles are obtained.

The fifth section deals with certain transfer characteristics giving valuable response data in special cases.

It should be noted again that the presentation used in this chapter emphasizes the essential unity of presently existing methods of servo analysis, and such methods are considered directly in terms of transfer function graphical representations. This basic unity is stressed throughout the chapter, and concepts such as frequency response phase margin and gain margin, which are frequently used in the literature, are mentioned only incidentally as items of interest in special cases.

SECTION 2 - THE GENERALIZED NYQUIST METHOD

The first, and most approximate, method to be discussed utilizes polar transfer function plots together with the closed-loop s-plane plot of transfer function poles and zeros. By using the results of a mapping theorem it is possible to consider a region of the closed-loop s-plane and determine the number of poles of the closed-loop transfer function within that region.

Before this method can be developed some fundamental concepts must be understood. The first concept is that of the "closed-loop s-plane." The second involves a basic mapping theorem. These ideas will be discussed initially, followed by a development and application of the method discussed above to the problem of determining stability.

(a) CLOSED LOOP S-PLANE.

As pointed out in chapter II the poles and zeros of any transfer function may be plotted on an s-plane. Such a plot may then be considered a graphical representation of the transfer function. If the poles and zeros of a closed-loop transfer function are plotted on an s-plane, the plane has been particularized to the extent that it may now be referred to as a closed-loop s-plane. Furthermore the plot may be called a graphical representation of the closed-loop transfer function.

Since the problem at hand is to determine the closed-loop transfer function, and thus its poles and zeros, obviously the poles and zeros cannot be plotted explicitly on the closed-loop s-plane. However, it is known that the poles and zeros do exist.

It is the aim of the balance of this section to develop a technique by which the poles and zeros of the closed-loop transfer function may be located approximately on the closed-loop s-plane.

(b) THE MAPPING THEOREM.

As pointed out previously, in order to determine the behavior of a closed-loop system, it is necessary to locate in the s-plane the poles and zeros of the expression $Y/(1+Y)$.

One method, to be described here, of locating these zeros and poles requires the use of two graphical constructions. These constructions will be illustrated by example.

Consider the simple third order system of (III-4):

$$(III-4) \quad Y = \frac{K}{s(as+1)(bs+1)}$$

Y may be plotted for values of s described by the contour shown in figure III-2. (The choice of the s contour may be considered arbitrary in this example,

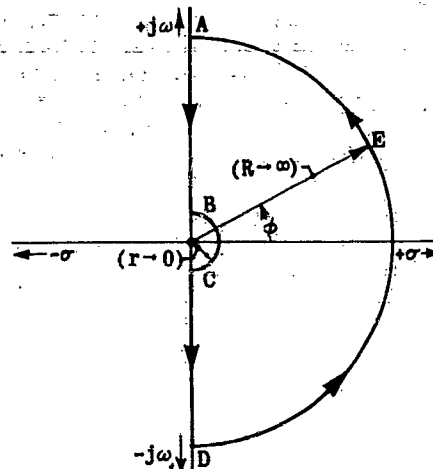


Figure III-2. S-Plane Contour

but the reasons for choosing particular contours in the s-plane will be explained later.) Proceeding from A to B, $s = +j\omega$ and $Y(s)$ appears as the heavy solid part of figure III-3.

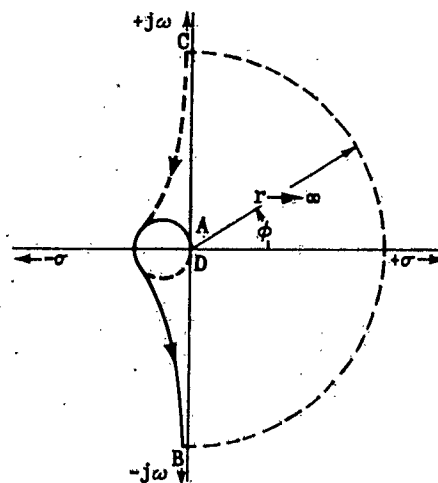


Figure III-3. $Y(s)$ -Plane Mapping

Similarly in the region C to D of figure III-2 $s = -j\omega$ and $Y(s)$ appears as the heavy dotted trace of figure III-3*. The segment B to C on the s -plane contour is a semi-circle of small radius ($r \rightarrow 0$) designed to exclude the pole at the origin ($s = 0$). To determine how this plots in the $Y(s)$ -plane substitute $s = re^{j\phi}$ in (III-4) then

$$Y(re^{j\phi}) = \frac{K}{re^{j\phi}(are^{j\phi} + 1)(bre^{j\phi} + 1)}$$

Hence:

$$(III-5) \quad \lim_{r \rightarrow 0} Y(re^{j\phi}) = \infty e^{-j\phi}$$

From (III-5) it may be seen that, as r swings from B ($\phi = +\pi/2$) to C ($\phi = -\pi/2$), Y swings from $-\infty$ to $+\infty$ in a positive sense. This is shown by the light dotted arc in figure III-3. The remainder of the s -plane contour (arc DEA), figure III-2, represents $\lim_{r \rightarrow \infty} Y(re^{j\phi})$ and maps into the origin on the $Y(s)$ -plane (III-3).

$$(III-6) \quad \lim_{R \rightarrow \infty} Y(Re^{j\phi}) = 0$$

The mapping theorem** states that: If the contour in the s -plane, figure III-2, positively encircles Z zeros and P poles of $1 + Y(s)$, the map of $Y(s)$ encircles the point $s = -1$ $N = Z - P$ times, where N is the number of encirclements and may be either positive or negative. A positive encirclement is defined as one in which the area enclosed by the contour is always on the left as the contour is traversed.

Since $1 + Y(s) = 1 + \frac{N(s)}{D(s)} = \frac{D(s) + N(s)}{D(s)}$, the poles (P) of $1 + Y(s)$ are simply the poles of $Y(s)$ and are known. Therefore, the zeros are determined from the equation:

$$(III-7) \quad Z = P + N$$

On figures III-2 and III-3 the arrows indicate the positive sense. The positive sense of encirclements can be remembered as the direction in which an observer would travel if he walked along the contour so that the interior was always to his left.

To apply these principles to the system of (III-4), first note that the contour of figure III-2 includes the entire right-half s -plane. Consequently, the examination of $Y(s)$ will determine the number of poles (P) in the region. To determine how many encirclements of the -1 point $Y(s)$ makes, draw a vector from the point -1 to the $Y(s)$ contour (figure III-4). As the arrow head moves along the contour, the vector pivots about its tail. Each time the vector sweeps out an angle of 2π in the positive (counter clockwise) direction, a positive encirclement is completed. In figure III-4 no encirclements occur ($N = 0$). Consequently, $Z = P + N = 0 + 0 = 0$.

Therefore, there are no roots of $1 + Y(s)$ in the right-half

* Note that the portion of the curve from C to D is the mirror image of that from B to A.

** For detailed proof of this theorem see the appendix to this volume, Section (A-IV).

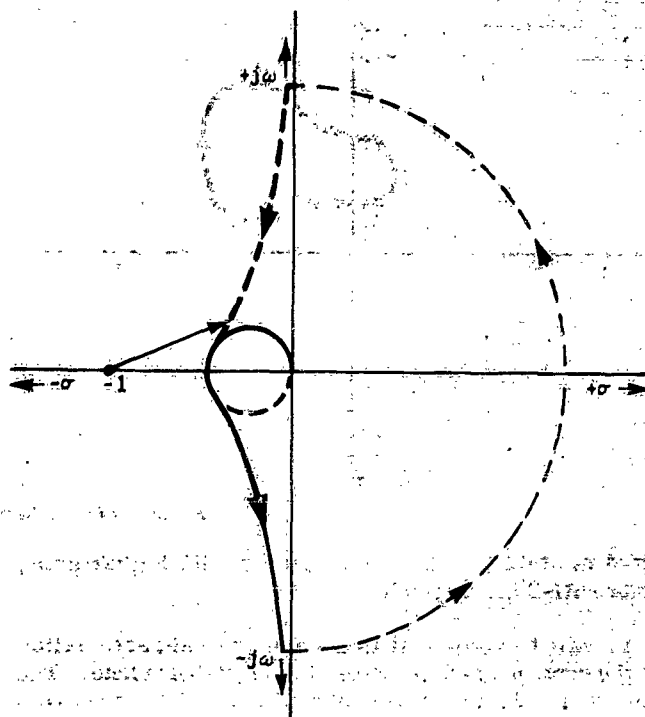


Figure III-4. Mapping

s -plane. Figure III-8b shows this same system plotted for a high enough value of the gain K so that there are two roots of $1 + Y(s)$ in the right-half of the s -plane. This is indicated by the encirclements of $s = -1$.

Although this theorem was discussed in terms of a specific problem, it is quite general. In the example chosen, only the right-half s -plane was mapped. But any region can be examined for the existence of zeros simply by properly designing the mapping contour (figure III-5).

Only two assumptions are necessary in order to apply the mapping theorem:

1. $Y(s)$ is a rational function of s .
2. None of the poles or zeros lies on the contour in the s -plane.

(It was because of assumption (2) that the contour in figure III-2 detoured around the pole at the origin.)

The following section discusses in detail the use of the mapping theorem to determine the stability of closed-loop systems.

(c) THE CONVENTIONAL NYQUIST CRITERION

It was pointed out in Chapter II that if any zeros of $1 + Y(s)$ have positive real parts, the system described by the open-loop transfer function $Y(s)$ is unstable. Hence, to verify stability alone, the entire right-half of the s -plane must be explored for zeros of $1 + Y(s)$ in the manner described in Section (b). (Now it is evident why the contour of figure III-2 was chosen in Section III-2b.)

Since there can be no zeros of $1 + Y(s)$ in this region if the system is to be stable, Z must be zero in $N = Z - P$ and hence the criterion for stability is that $N = -P$. The simple system described by the contour of figure

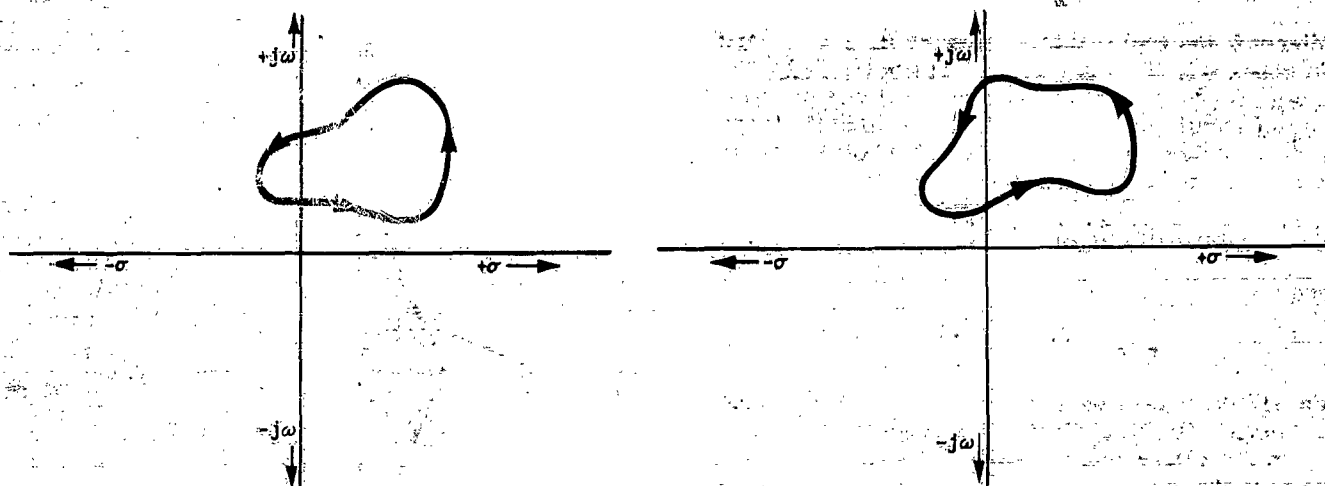


Figure III-5. General Contour Mapping

III-4 is stable. The same system with higher gain, figure III-8b, is unstable.

It is well to review at this point the characteristics of the contour used to enclose the right-half plane. The values of s defining this contour, figure III-2, are given as follows:

$$\begin{aligned}
 \text{(III-8)} \quad & \text{From A to B} \quad s = j\omega \quad \omega > 0^+ \\
 & \text{From B to C} \quad s = re^{j\phi} \quad r \rightarrow 0 \\
 & \quad \phi = \frac{\pi}{2} \text{ when } \omega = 0^+ \\
 & \quad \phi = -\frac{\pi}{2} \text{ when } \omega = 0^- \\
 & \text{From C to D} \quad s = -j\omega \quad 0^- > \omega > -\infty \\
 & \text{From D to E} \quad s = Re^{j\phi} \quad R \rightarrow \infty
 \end{aligned}$$

For the open-loop transfer function $Y(s)$, the presence of zeros of $1 + Y(s)$ in the right-half plane is readily determined by plotting $Y(s)$ with the above values substituted for s , and then simply counting the encirclements of -1 and then applying (III-7).

The actual plotting may be somewhat simplified by the following considerations: For physically realizable systems, the part of the s -plane contour DEA is unimportant in determining the plot of $Y(s)$. This is true since $Y(s) = \frac{N(s)}{D(s)} = \frac{A_n s^n + \dots + A_0}{\alpha_n s^n + \dots + \alpha_0}$

where n , the order of the denominator, must be greater than m , the order of the numerator. The value of $Y(s)$ corresponding to the contour from DEA is:

$$\begin{aligned}
 \text{(III-9)} \quad \lim_{R \rightarrow \infty} [Y(s)]_{s=Re^{j\phi}} &= \frac{A_n R^n e^{jn\phi} + \dots + A_0}{\alpha_n R^n e^{jn\phi} + \dots + \alpha_0} \\
 &= \lim_{R \rightarrow \infty} \frac{A_n}{\alpha_n R^{n-m}} = 0
 \end{aligned}$$

Consequently, the entire arc DEA maps into the origin of the $Y(s)$ plane.

The portion of the $Y(s)$ contour corresponding to the s -plane contour from A to B is just $Y(j\omega)$, the simplest form of transfer function plot developed in Chapter II. From C to D, $Y(s)$ is given by $Y(-j\omega)$, which is the mirror image about the real axis of $Y(+j\omega)$.

The only part of the contour remaining is that from B to C. In many engineering problems, Y has a pole of order n at the origin. But one of the conditions upon which the mapping theorem can be applied is that the s -plane contour does not pass through any poles. Consequently, the contour is detoured by letting $s = re^{j\phi}$ (with r very small) near the origin.

Then, since $Y(s)$ is of the form

$$\text{(III-10)} \quad Y(s) = \frac{N(s)}{s^n D(s)}$$

substituting $s = re^{j\phi}$ to avoid the pole of order n

$$\text{(III-11)} \quad Y(s) = \lim_{r \rightarrow 0} Y(re^{j\phi}) = \lim_{r \rightarrow 0} \frac{A_n r^n e^{jn\phi} + \dots + A_0}{r^n e^{jn\phi} [\alpha_n r^k e^{jk\phi} + \dots + \alpha_0]}$$

Now as the s -plane contour in the region B to C is traversed in a positive direction, ϕ goes from $+\pi/2$

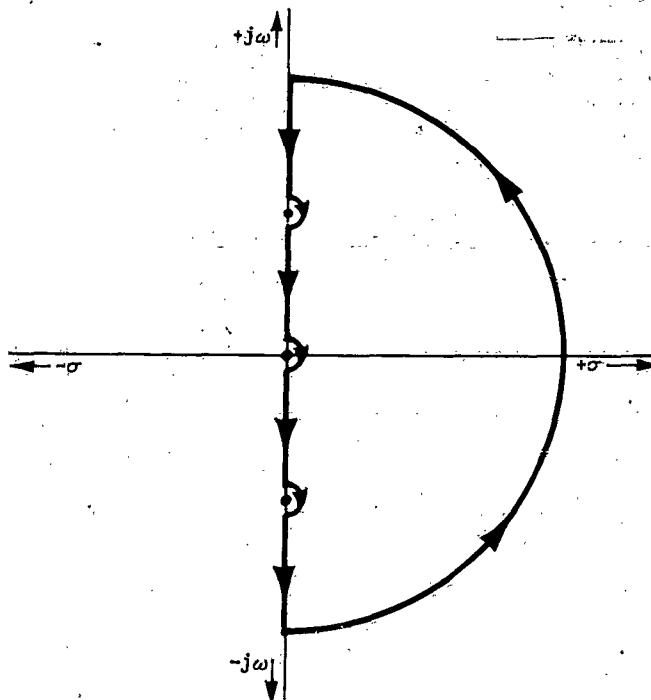


Figure III-6. Poles on Imaginary Axis

to $-\pi/2$. Consequently, the contour of $Y(s)$ starts at $-\pi/2$ and proceeds to $+\pi/2$, in a positive sense and at a radius equal to $\lim_{r \rightarrow 0} \frac{A_0}{\alpha_0 r^n} = \infty$. Evidently,

the portion of $Y(s)$ corresponding to the one from B to C of the s -plane contour is an arc of infinite radius.

There are two common situations in which this procedure must be slightly modified. First, if A_0/α_0 is negative, the $Y(s)$ contour sweeps from $+\pi/2$ to $-\pi/2$ in a positive sense as s traverses the s -plane contour from B to C, figure III-2. Secondly, it sometimes occurs that the denominator of $Y(s)$ includes factors of the form $s^2 + \omega^2$. To avoid these poles and those at the origin the s -plane contour must take the shape of figure III-6 and the map of $Y(s)$ be treated accordingly.

If $Y(s)$ has no poles at the origin, the device of letting $s = re^{j\phi}$ near $s = 0$ is not required, and the entire contour of $Y(s)$ is made up of $Y(+j\omega)$ and $Y(-j\omega)$, figure III-7.

The above points are illustrated in figures III-8a to III-8g.

(d) SPECIFIED MINIMUM DAMPING AND DAMPING RATIO.

It is clear that the mapping theorem discussed in section (b) of this chapter could be used in checking any desired region of the s -plane for closed loop poles. However,

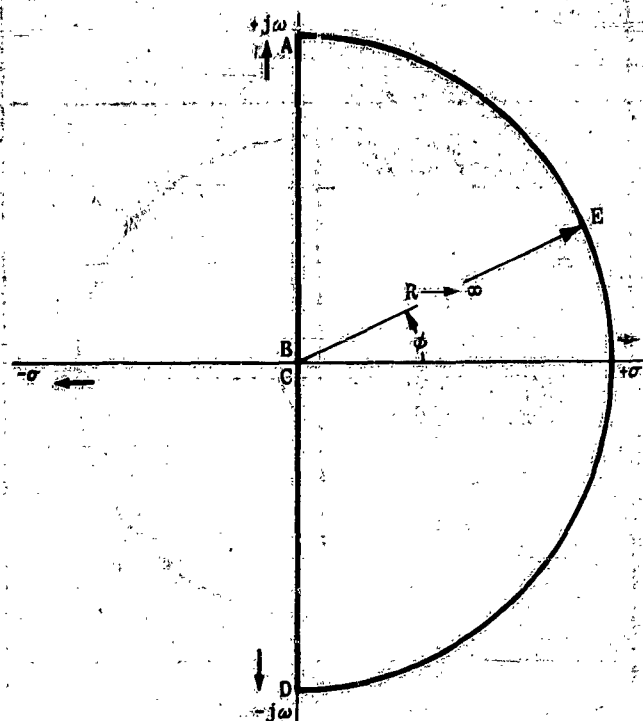


Figure III-7. Simple Contour

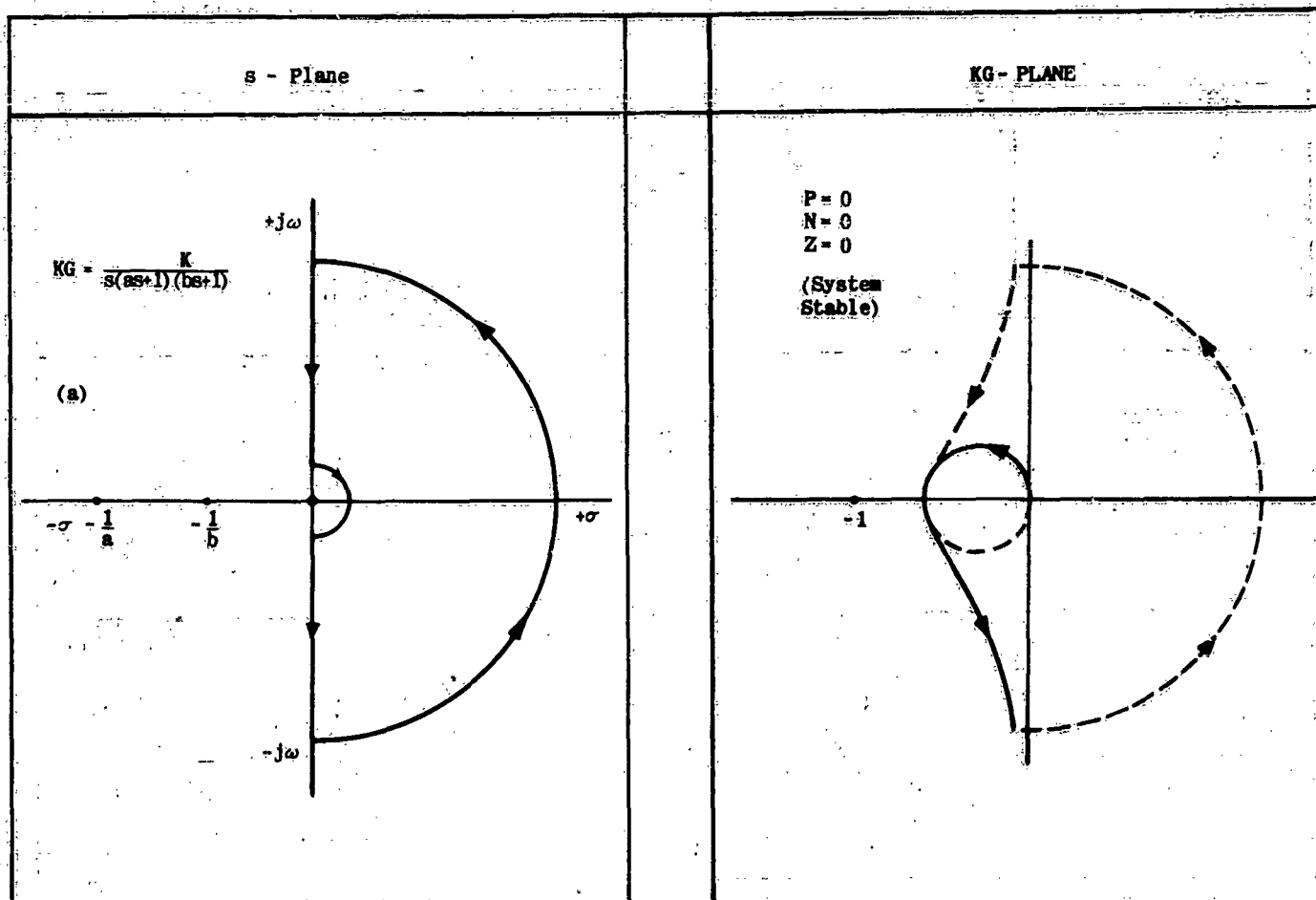


Figure III-8 (Sheet 1 of 3 Sheets). Examples of Mapping for Stable and Unstable Systems

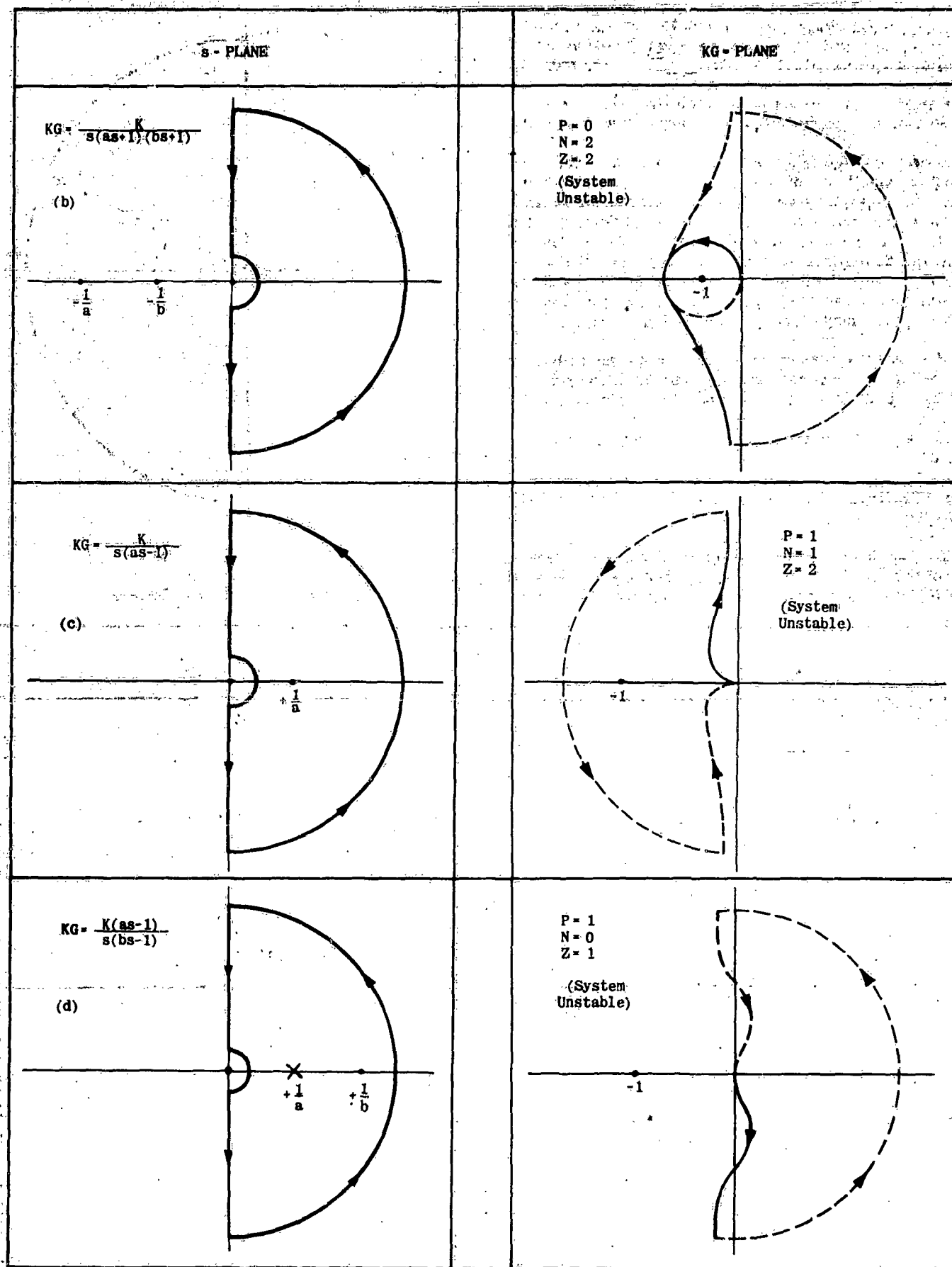


Figure III-8 (Sheet 2 of 3 Sheets). Examples of Mapping for Stable and Unstable Systems

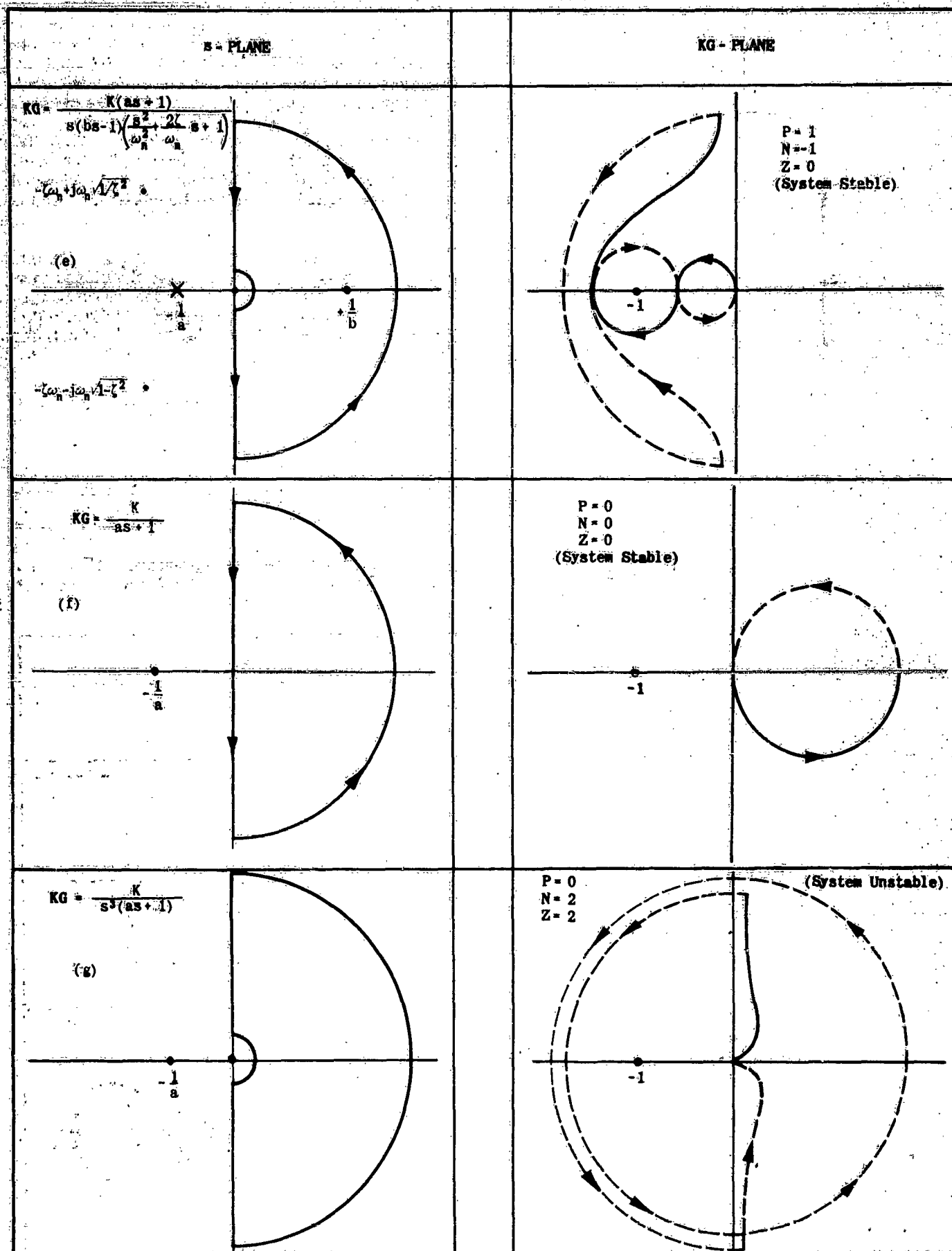


Figure III-8 (Sheet 3 of 3 Sheets). Examples of Mapping for Stable and Unstable Systems

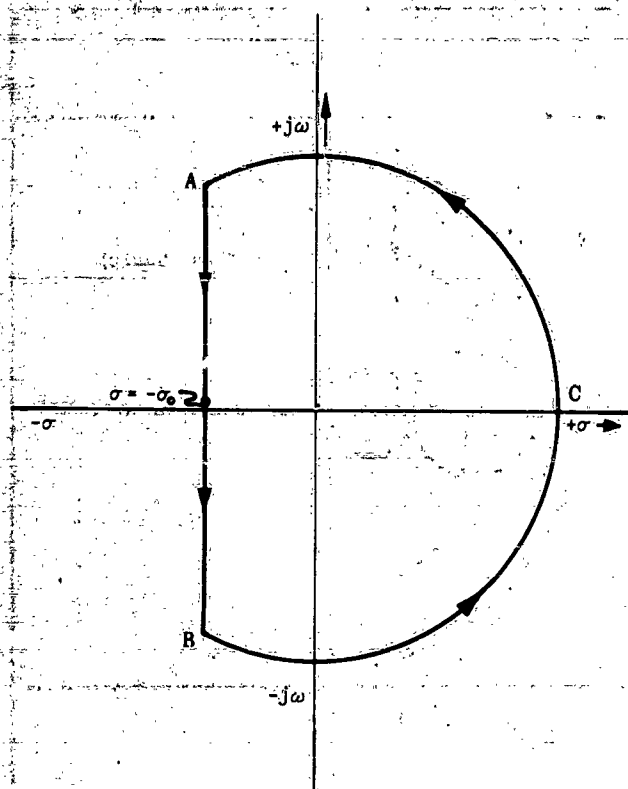


Figure III-9. Minimum Damping Contour

in all cases, the regions checked define the values of s to be substituted in $Y(s)$, and hence practical use of analysis time usually requires that the boundaries of the regions in the s -plane to be examined be defined by comparatively simple functions of s . The simplest case, that of the conventional Nyquist criterion covered in section III-2c, is of much value in checking stability. Two other cases, those of minimum damping and minimum damping ratio are also of interest and occasional importance. These cases will be considered in this section.

First, consider the minimum damping case. In this situation it is desired to investigate the closed loop transfer function for poles having damping less than some specified value. The s -plane contour in this case is shown in figure III-9. The values of s defining the contour are:

$$(III-12) \quad \text{From A to B} \quad s = -\sigma_0 + j\omega \quad \omega \geq \omega \geq -\omega$$

$$\text{From B to A} \quad s = R e^{j\phi} \quad R \rightarrow \infty$$

where $1/\sigma_0$ = min. damping time constant $\phi = -\pi/2$ when $\omega = -\infty$
 $\phi = +\pi/2$ when $\omega = +\infty$

For the same reasons mentioned previously, the portion of the contour BCA does not contribute to the locus of

s - Plane	KG (s)
$KG = \frac{K}{ab} \frac{1}{s(s+\frac{1}{a})(s+\frac{1}{b})}$	
	<p>If $(-1+j0)$ lies between A and B $P = 3, N = -1, Z = 2$ If $(-1+j0)$ lies between A and C $P = 3, N = 0, Z = 3$</p> <p>(a)</p>
	<p>If $(-1+j0)$ lies between A and B $P = 2, N = 0, Z = 2$</p> <p>(b)</p>

Figure III-10 (Sheet 1 of 2 Sheets). Examples of Minimum Damping

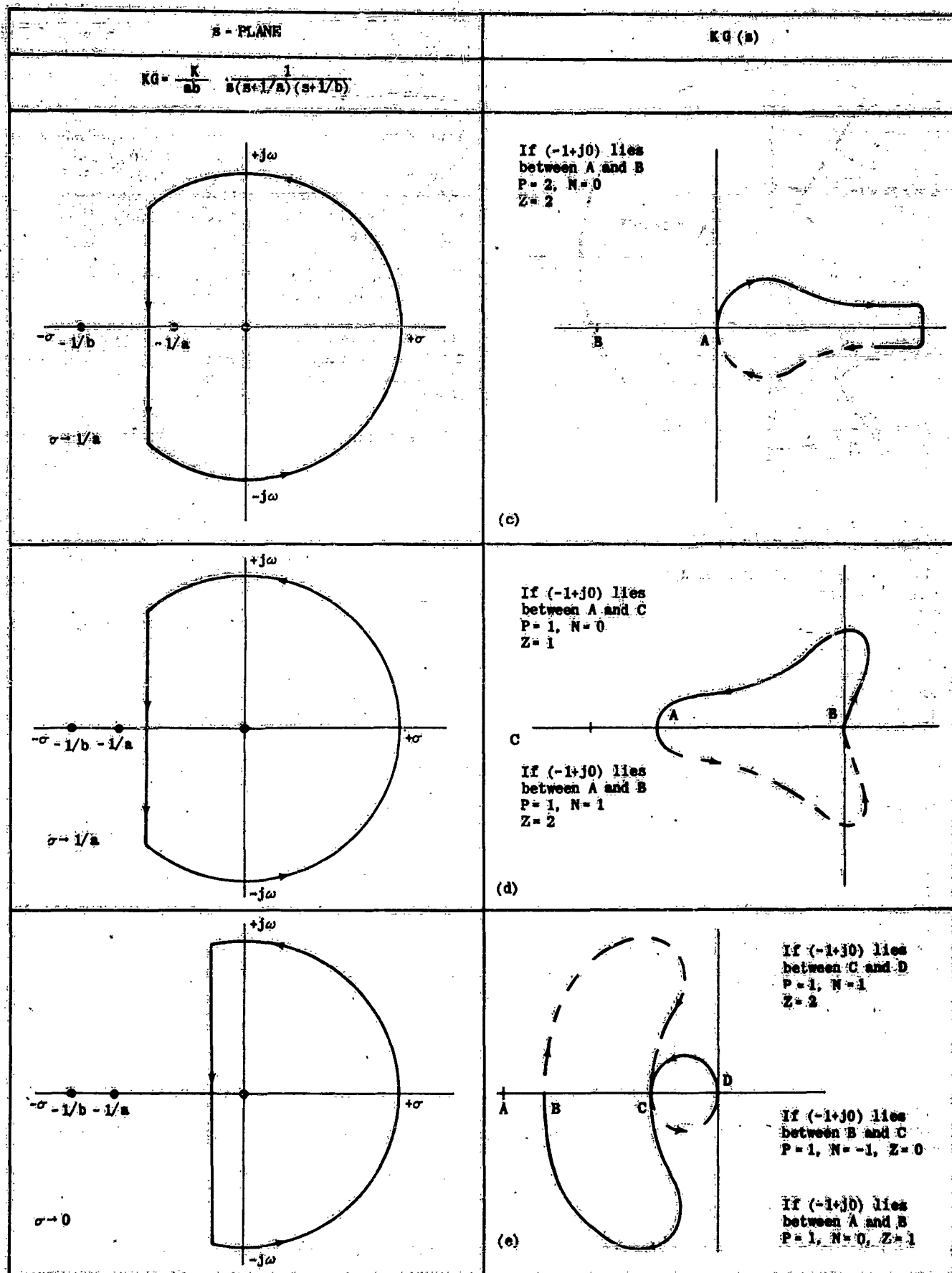


Figure III-10 (Sheet 2 of 2 Sheets). Examples of Minimum Damping

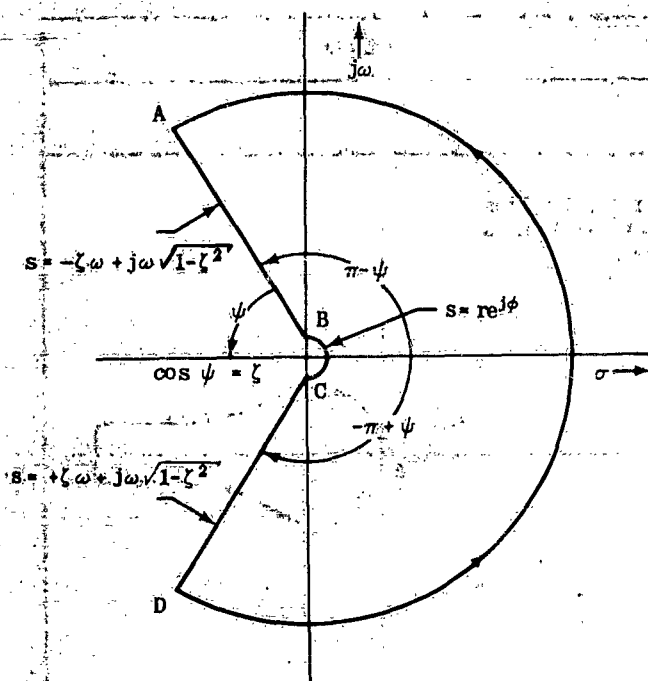


Figure III-11. Minimum Damping Ratio Contour

$Y(s)$ for physically realizable systems. The value of $Y(s)$ corresponding to the remainder of the contour in the s -plane is given by:

$$(III-13) \quad Y(s) = Y(-\sigma_0 + j\omega)$$

Typical examples for this situation are shown in figure (III-10).

SECTION 3—THE OPEN LOOP—CLOSED LOOP LOGARITHMIC METHOD

(a). GENERAL

This section discusses a method of estimating closed loop poles and zeros from logarithmic open loop transfer function plots.

If the transfer functions concerned are ratios of rational polynomials, they may be written in factored form, all factors in either numerator or denominator being first order or quadratic. Each of the magnitudes (absolute values) of these factors may be represented by a pair of asymptotes on a logarithmic plot, except near the intersection of these asymptotes (the "break-point"). Then, since the use of logarithmic coordinates permits addition of the logarithms to replace multiplication (and division) of the factors, by adding the asymptotic factor plots graphically, an "asymptotic" graph of the entire transfer function results. All the straight line approximations obtained in this way will be referred to in this volume as asymptotes. It is realized that this is an extension of the strict mathematical meaning of the word asymptote, but this practice is well justified by usage in the controls field.

It will be recalled that the intersection of two asymptotes with a slope difference of 20 db/dec indicates the presence of a first order pole or zero, with the direction

Note that for damping greater than the minimum chosen, there must be no zeros of $1+Y(s)$ in the domain enclosed by the contour in the s -plane; i.e., in the $Y(s)$ plane the number of encirclements of the -1 point must be such that $Z=N+P=0$.

It is interesting to note in the examples given that under no circumstances can the minimum damping time constant be less than a , (see figure III-10e), and that to obtain even this damping time the gain must be adjusted so that the -1 point lies between points B and C in figure III-10e.

For the specified minimum damping ratio case the s -plane contour is shown in figure III-11.

Note that possible poles of $1+Y(s)$ at $s=0$ (and elsewhere on the contour) must be avoided in this case as in the conventional Nyquist case. The values of s defined by the contour are given by

$$(III-14) \quad \text{From A to B} \quad s = -\zeta\omega + j\omega\sqrt{1-\zeta^2} \quad \omega > 0^+$$

$$\text{From B to C} \quad s = re^{j\phi} \quad r \rightarrow 0$$

$$\phi = \pi - \psi \quad \text{when} \quad \omega = 0^+$$

$$\phi = -\pi + \psi \quad \text{when} \quad \omega = 0^-$$

$$\text{From C to D} \quad s = \zeta\omega + j\omega\sqrt{1-\zeta^2} \quad 0^- > \omega > -\infty$$

$$\text{From D to A} \quad s = Re^{j\phi} \quad R \rightarrow \infty$$

$$\text{and} \quad \phi = -\pi + \psi \quad \text{when} \quad \omega = -\infty$$

$$\phi = \pi - \psi \quad \text{when} \quad \omega = +\infty$$

A typical application of the minimum damping ratio case is shown in figure III-12.

of the phase curve in that region indicating whether the pole or zero is of minimum or non-minimum phase (i.e., in the left or right half of the s plane representation). For quadratic factors (slope differences at intersection of 40 db/dec), the asymptote intersection occurs at the undamped natural frequency, and the departure of the actual plot from the asymptote intersection is equal to twice the damping ratio in decibels. These characteristics of logarithmic transfer function plots are basic to this section. Broadly speaking, the method involves the following steps:

1. The open-loop transfer function, $Y(s)$, is plotted logarithmically with s set equal to $j\omega$.
2. For a given open-loop gain, K , (or open-loop zero db line), the closed loop transfer function plot is constructed. This construction is materially aided by the use of approximations and a special chart.
3. The analytical factored form of the closed loop transfer function is obtained by establishing the asymptotic representation of the plotted curve, and by utilizing the characteristics of the logarithmic plot. This process is aided by the use of previously known data (the zeros of the closed loop transfer function).

The first item of the above list has been extensively

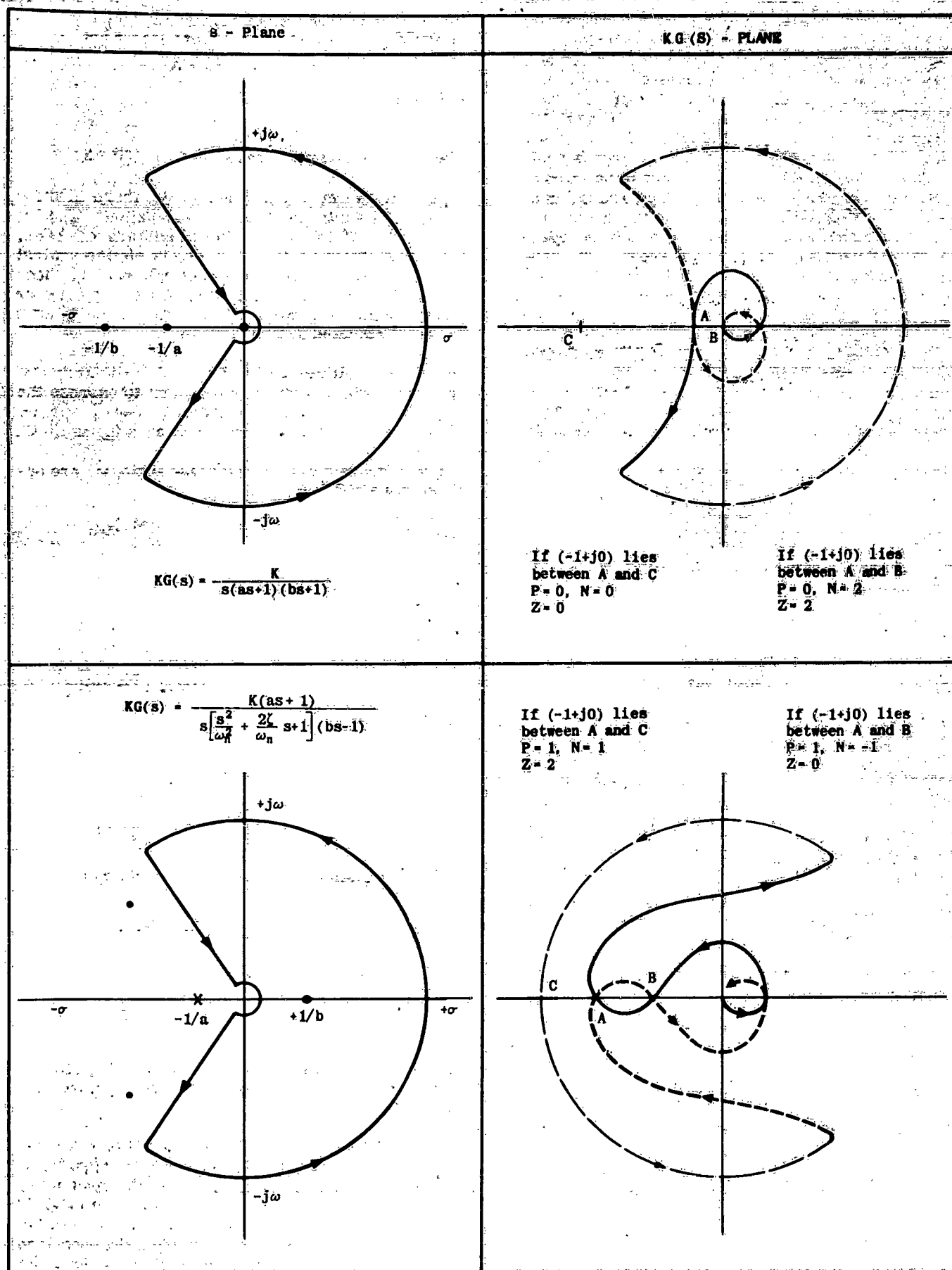


Figure III-12. Examples of Minimum Damping Ratio

Chapter III

Section 3

discussed in the transfer function section of chapter II, and needs no further explanation. The second item can be achieved by use of approximations and the use of a special chart where required. A subsection will be devoted to the derivation of this chart and approximations with examples of their use. The third item (finding the analytic form of a given closed loop plot) is essentially a curve fitting task. The curve fitting process is made practical by the fact that the asymptotes to be fitted to the closed loop plot have slopes which are integral multiples of ± 20 db/dec. Another aid is the fact that the zeros of the closed loop are also the zeros of the open-loop transfer function, and hence are known. The process of finding an analytic factored form of a given transfer function plot is discussed, with examples, in a further subsection.

(b) RELATIONSHIPS BETWEEN OPEN AND CLOSED-LOOP TRANSFER FUNCTIONS—CHARTS

The portion of the closed loop transfer function, $Z(s)$, of interest to the discussion of this chapter is given by

$$(III-15) \quad Z(s) = \frac{Y(s)}{1 + Y(s)}$$

A method will now be developed to relate $Y(s)/[1+Y(s)]$ to $Y(s)$ by simple graphical means. For convenience $Y(j\omega)/[1+Y(j\omega)]$ will be related to $Y(j\omega)$. Since this is merely a functional representation, $j\omega$ may later be replaced by s .

$Y(j\omega)$ is usually readily available in the logarithmic graphical form. The logarithmic graphical representation of $Y(j\omega)/[1+Y(j\omega)]$ is obtained from $Y(j\omega)$ by using a special chart which will now be derived.

If the magnitude (amplitude ratio) and the phase angle of the closed-loop transfer function are denoted by M and ψ , respectively, relationships between $Y(j\omega)$ and constant values of M and ψ can be derived. If, in addition, contours of constant values of M and ψ are superimposed upon a linear, rectangular plot of $Y(j\omega)$, amplitude ratio of $Y(j\omega)$ plotted versus the phase angle of $Y(j\omega)$, the phase angle and amplitude ratio of $Y(j\omega)/[1+Y(j\omega)]$ can be read directly off the plot for a given value of frequency ω . In the following paragraphs these contours and their associated chart will be developed.

The open-loop transfer function, $Y(j\omega)$, may be represented as a complex number

$$(III-16) \quad Y(j\omega) = x(\omega) + jy(\omega)$$

Then

$$(III-17) \quad M = \left| \frac{Y(j\omega)}{1+Y(j\omega)} \right| = \left| \frac{x(\omega) + jy(\omega)}{1+x(\omega) + jy(\omega)} \right|$$

$$= \left\{ \frac{[x(\omega)]^2 + [y(\omega)]^2}{[1+x(\omega)]^2 + [y(\omega)]^2} \right\}^{1/2}$$

$$M^2 \{ [1+x(\omega)]^2 + [y(\omega)]^2 \} = [x(\omega)]^2 + [y(\omega)]^2$$

$$[y(\omega)]^2 (M^2 - 1) + [x(\omega)]^2 (M^2 - 1) + [x(\omega)]^2 2M^2 + M^2 = 0$$

By completing the square of the x terms:

$$(III-18) \quad [y(\omega)]^2 + \left[x(\omega) + \frac{M^2}{M^2-1} \right]^2 = \frac{M^2}{(M^2-1)^2}$$

(III-18), plotted on a linear rectangular plot, y vs. x , is that of a series of circles, with centers at

$$x_0 = -\frac{M^2}{M^2-1}, \quad y_0 = 0$$

and with radii

$$r = \left| \frac{M}{M-1} \right|$$

These "M circles" are shown in figure III-13.

Note that to any fixed value of M , there corresponds an infinite number of combinations of $x(\omega)$ and $y(\omega)$. If the log of the magnitude of $Y(s)$, i.e., $\log [x^2(\omega) + y^2(\omega)]^{1/2}$, in db, is plotted against $\tan^{-1}[y(\omega)/x(\omega)]$ for a series of values of M , these "M circles" or lines of constant closed loop transfer function amplitude ratio, plot into "M contours" such as those indicated in figure III-14.

Since the amplitude ratio of $Y(j\omega)$ is ordinarily expressed in db, it is most convenient to express the amplitude ratio of $Z(j\omega)$ in this form. Consequently on all the ensuing plots, M is given in db as in figure III-15.

Lines of constant closed loop phase angle, ψ , are derived in a similar fashion.

$$(III-19) \quad \angle \left[\frac{Y(j\omega)}{1+Y(j\omega)} \right] = \psi = \tan^{-1} \frac{y(\omega)}{x(\omega)} - \tan^{-1} \frac{y(\omega)}{1+x(\omega)}$$

By trigonometry, then:

$$(III-20) \quad \psi = \tan^{-1} \left\{ \frac{\frac{y(\omega)}{x(\omega)} - \frac{y(\omega)}{1+x(\omega)}}{1 + \left[\frac{y(\omega)}{x(\omega)} \right] \left[\frac{y(\omega)}{1+x(\omega)} \right]} \right\}$$

which simplifies to

$$(III-21) \quad \tan \psi = \frac{y}{x^2 + x + y^2}$$

$$\text{or} \quad (x + \frac{1}{2})^2 + \left(y - \frac{1}{2 \tan \psi} \right)^2 = \frac{1}{4} \left[\frac{\tan^2 \psi + 1}{\tan^2 \psi} \right]$$

again the equation of a circle. The radius is

$$(III-22) \quad r = \frac{1}{2 \tan \psi} \sqrt{\tan^2 \psi + 1}$$

and the center is located at

$$(III-23) \quad x_0 = -1/2; \quad y_0 = \frac{1}{2 \tan \psi}$$

These circles, when plotted on the same logarithmic coordinates as the M contours, are the ψ contours of figure III-14.

The plot of figure III-14, commonly called a Nichols chart, enables the analyst to find easily the closed-loop values of the transfer function from the open-loop, either by actual plotting on the chart, or simply by reference. An example of this method is shown in figure III-15, going successively from the open-loop plot shown in figure III-15a, to the Nichols chart plot of figure III-15b, and then to the closed loop plot of figure III-15c.

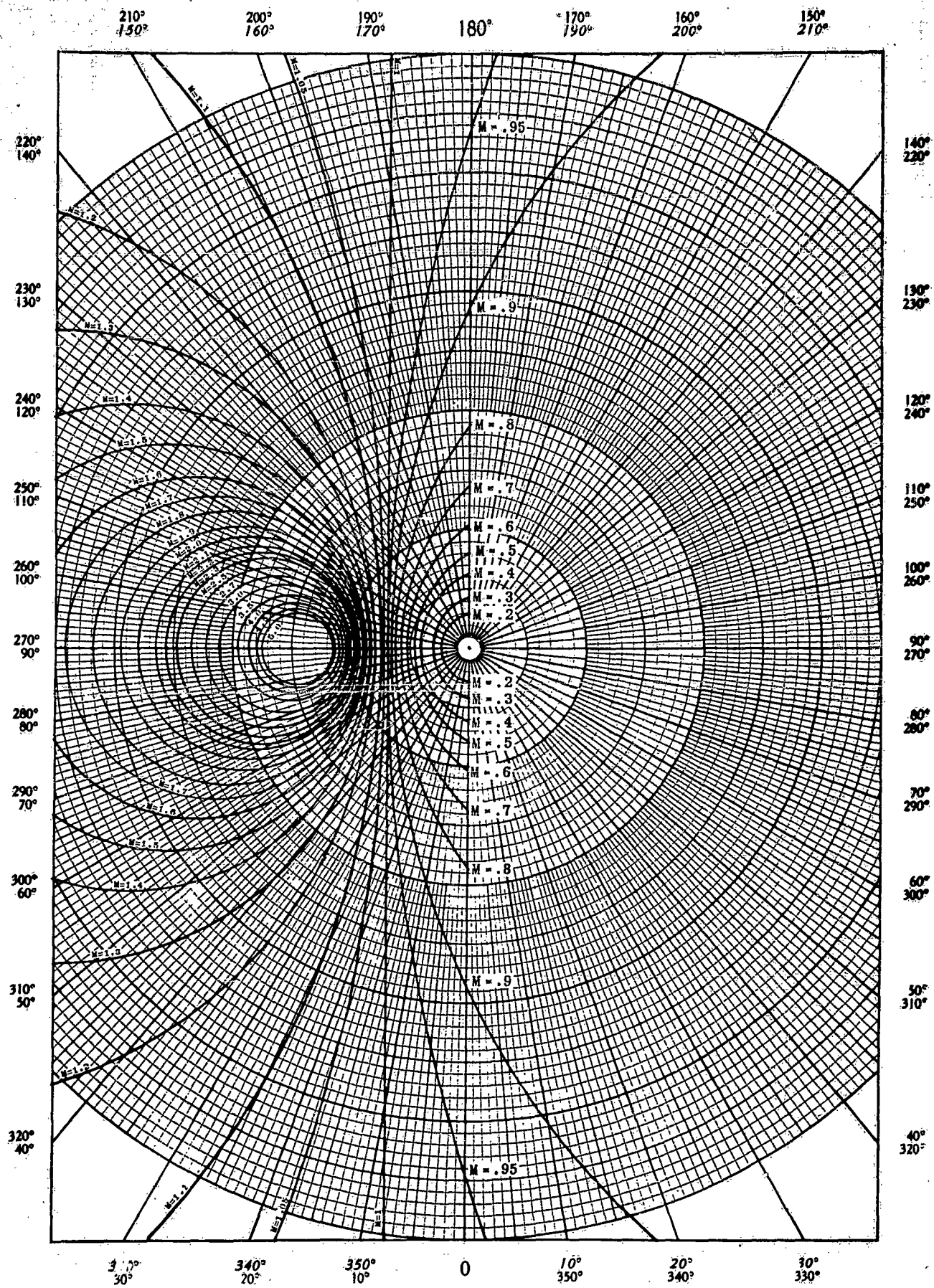


Figure III-13. M - Circles

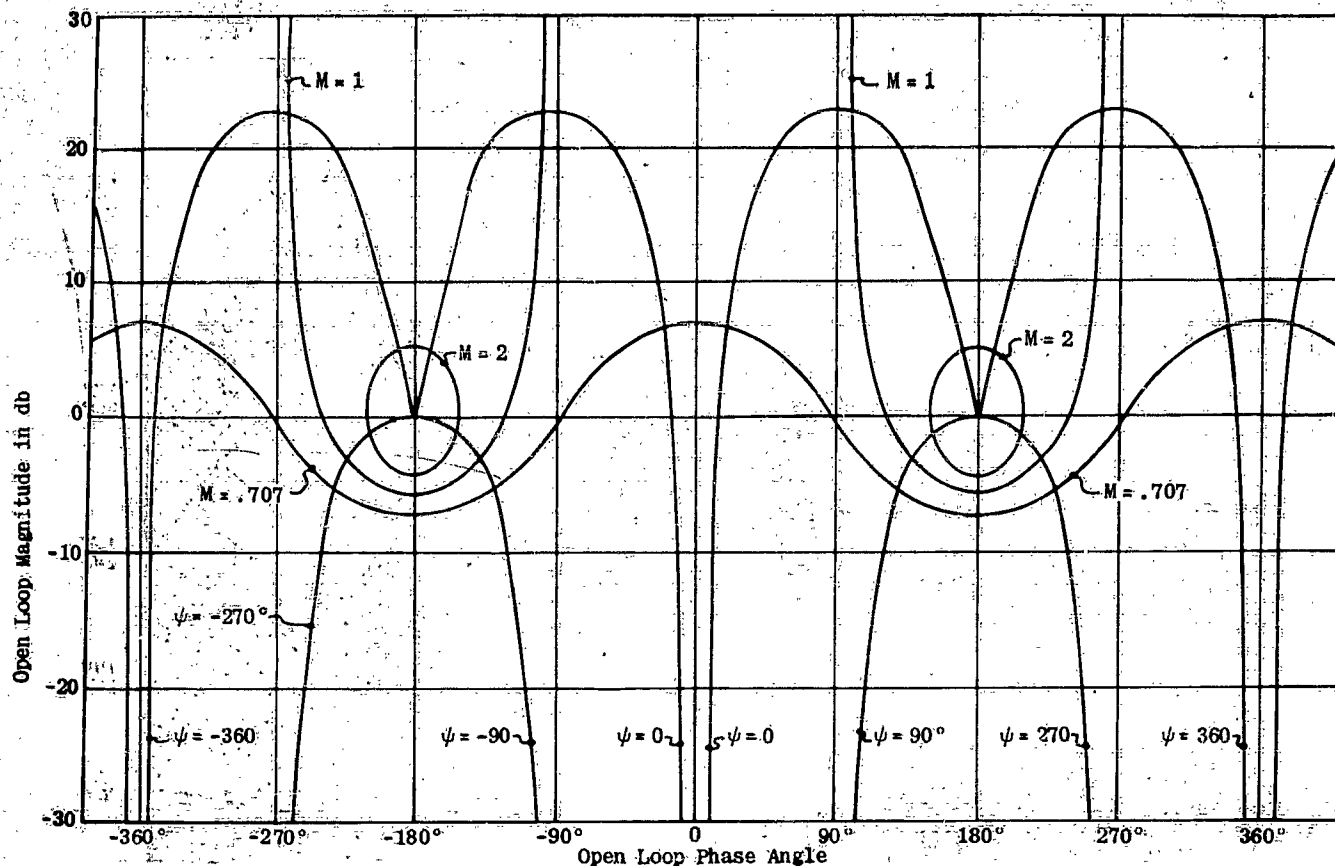


Figure III-14. Nichols Chart

In practice, the actual plotting shown in figure III-15b may be avoided by utilizing the chart as a graphical table.

Accurate versions of Nichols charts are provided in figures A-18 and A-19 of the appendix.

(c) RELATIONSHIP BETWEEN OPEN AND CLOSED-LOOP TRANSFER FUNCTIONS—BY APPROXIMATION

Because of certain approximate relationships, much of the work required to obtain $Y(j\omega)/[1+Y(j\omega)]$ from $Y(j\omega)$ by the method shown in the previous subsection (III-3b) may be reduced. These approximate relationships are especially valuable in preliminary studies and for any case where rapid, though not extremely accurate results are required.

It is evident from (III-15) that if $|Y(s)| \gg 1$ then

$$(III-24) \quad |Z(s)| \approx 1$$

or zero decibels. Also if $|Y(s)| \ll 1$

$$(III-25) \quad |Z(s)| \approx |Y(s)|$$

An examination of a Nichols chart shows that if $|Y(j\omega)| \geq 25$ db, the relationship of (III-24) is correct within an error of 0.5db. If $|Y(j\omega)| \geq 10$ db, the error involved is of the order of 2 db.

However if $|Y(s)|$ is of the order of magnitude of 1

(0 db), the entire expression for $Z(s)$ must be used, i.e., where $|Y(s)| \sim 1$

$$(III-26) \quad Z(s) = Y(s)/[1 + Y(s)]$$

In those regions of $Y(s)$ where (III-24) or (III-25) apply, the plot of $Z(s)$ is given directly by a knowledge of $Y(s)$. In those regions where (III-26) applies, the method outlined in section III-3b must be used.

For illustrative purposes several examples of obtaining closed-loop plots from open-loop plots are given below.

EXAMPLE 1. Assume an open-loop transfer function $3.16/(j\omega + 1)$ which is drawn lightly in figure III-16. Where $|Y(j\omega)| = |3.16/(j\omega + 1)| \geq 10$ db it will be assumed that the closed-loop transfer function $|Z(j\omega)| = 0$ db. Where $|Y(j\omega)| \leq -10$ db it will be assumed that $|Z(j\omega)| = |Y(j\omega)|$. Thus, the dark lines are drawn in to represent $|Z(j\omega)|$ in these regions. Where $10 \text{ db} > |Y(j\omega)| > -10 \text{ db}$ the Nichols chart has been used, and the results are indicated by dashed lines. It is evident that the dashed line does not fair into the solid line at the lower end of the frequency range. This is a result of the approximation that was made above and could have been predicted. This difficulty may be overcome at the outset by a simple calculation of the steady state gain. Since

$$Y(s) = K \frac{N(s)}{D(s)}, \quad Z(s) = \frac{KN(s)}{D(s) + KN(s)}$$

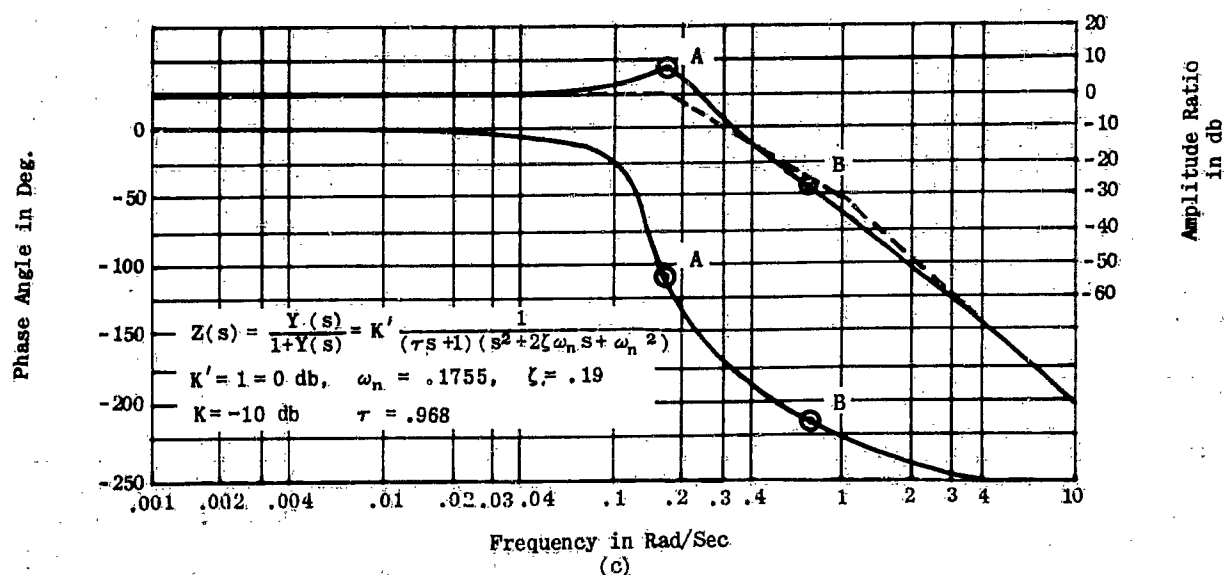
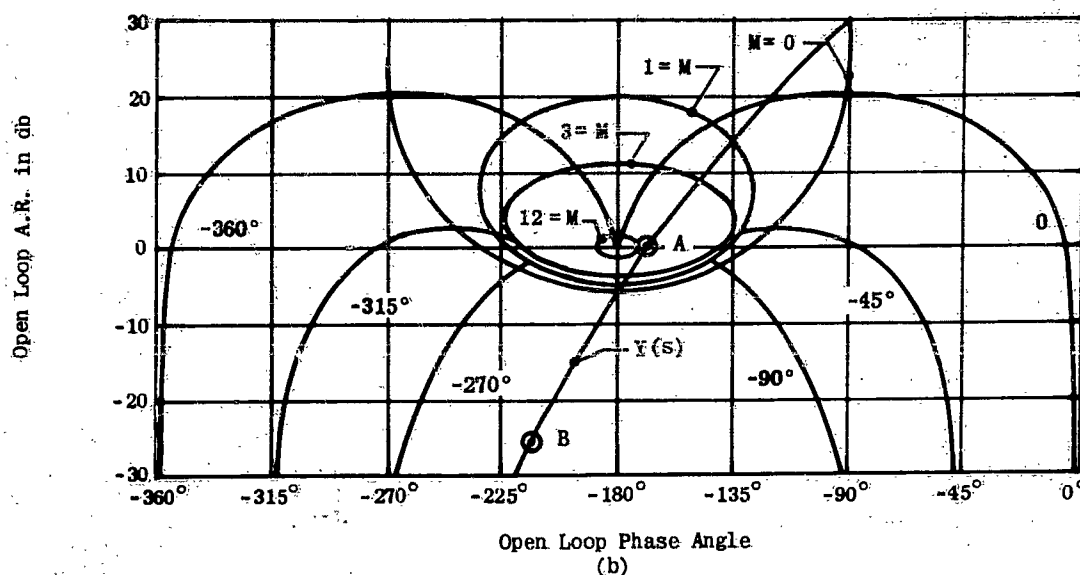
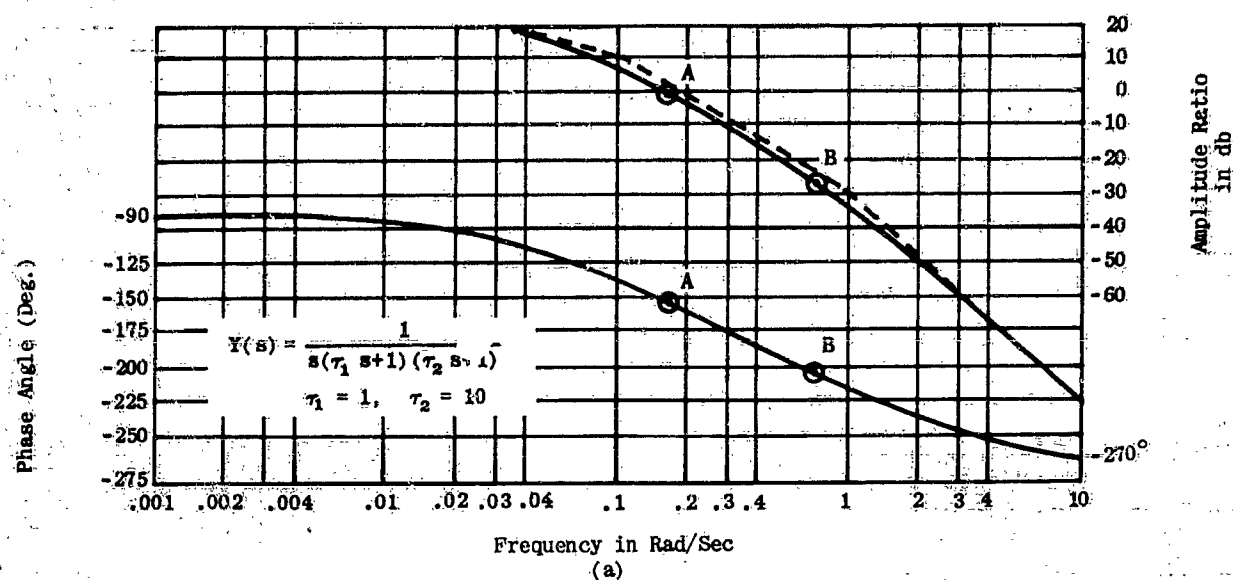


Figure III-15. Open Loop to Closed Loop by Nichols Chart

finer by the breakpoints of an asymptotic plot if all the poles and zeros are first order. If a pole or zero is of second order, the departure of the actual curve from the breakpoint is also required. The sign of the pole or zero, i.e., whether the term is minimum or non-minimum phase, is completely determined by the direction or slope of the phase curve in the region of the pole or zero. Therefore, to determine the analytic form of the closed-loop function from its plot, it is only necessary to establish the asymptotic plot, and to observe the phase angle change in the immediate vicinity of break points. Values of poles and zeros are then determined directly from the break point frequencies (and break point departures in the case of quadratic terms), and signs are determined by the local phase angle change. The functional form of $Y(j\omega)/[1+Y(j\omega)]$ is then known analytically, and $Y(s)/[1+Y(s)]$ is obtained simply by substituting s for $j\omega$.

In determining the asymptotic plot, the following facts are of value:

1. Any asymptote is a straight line with a slope which is an integral multiple of ± 20 db/dec. Therefore, for any transfer function expressible as the ratio of rational, constant coefficient polynomials, the change in slope of the asymptotic plot at any break point must be an integral multiple of ± 20 db/dec. The value of the integral multiplier is, of course, the order of the pole or zero.

2. In regions where equations (III-24) and (III-25) are valid, a portion of the approximate analytic factored form of the closed-loop is known by a simple inspection of the open-loop logarithmic transfer function plot.

3. The zeros of the closed-loop transfer function or its analytic form are known initially, since they are also the zeros of the open-loop transfer function.

To illustrate the methods of obtaining the closed-loop transfer function in analytic, factored form several examples are given below.

EXAMPLE 1. A closed loop transfer function is shown in figure III-19. If the asymptotes to the amplitude curve are drawn, it is found that, when extended, they meet at a frequency $1/\tau_1$. Also, at the frequency $1/\tau_1$ the phase is -45° . It is evident that the closed-loop transfer function has the form $Z(s) = K/(\tau_1 s + 1)$. From the amplitude curve it is seen that the gain in db is K' . The linear gain K may be obtained from figure A-20 which relates linear amplitude ratio to decibels. The complete closed-loop transfer function is then:

$$Z(s) = K/(\tau_1 s + 1)$$

EXAMPLE 2. This example is considerably more complex than the first. It is introduced to emphasize various details. The closed-loop transfer function is shown in figure III-20. The asymptotes to the amplitude

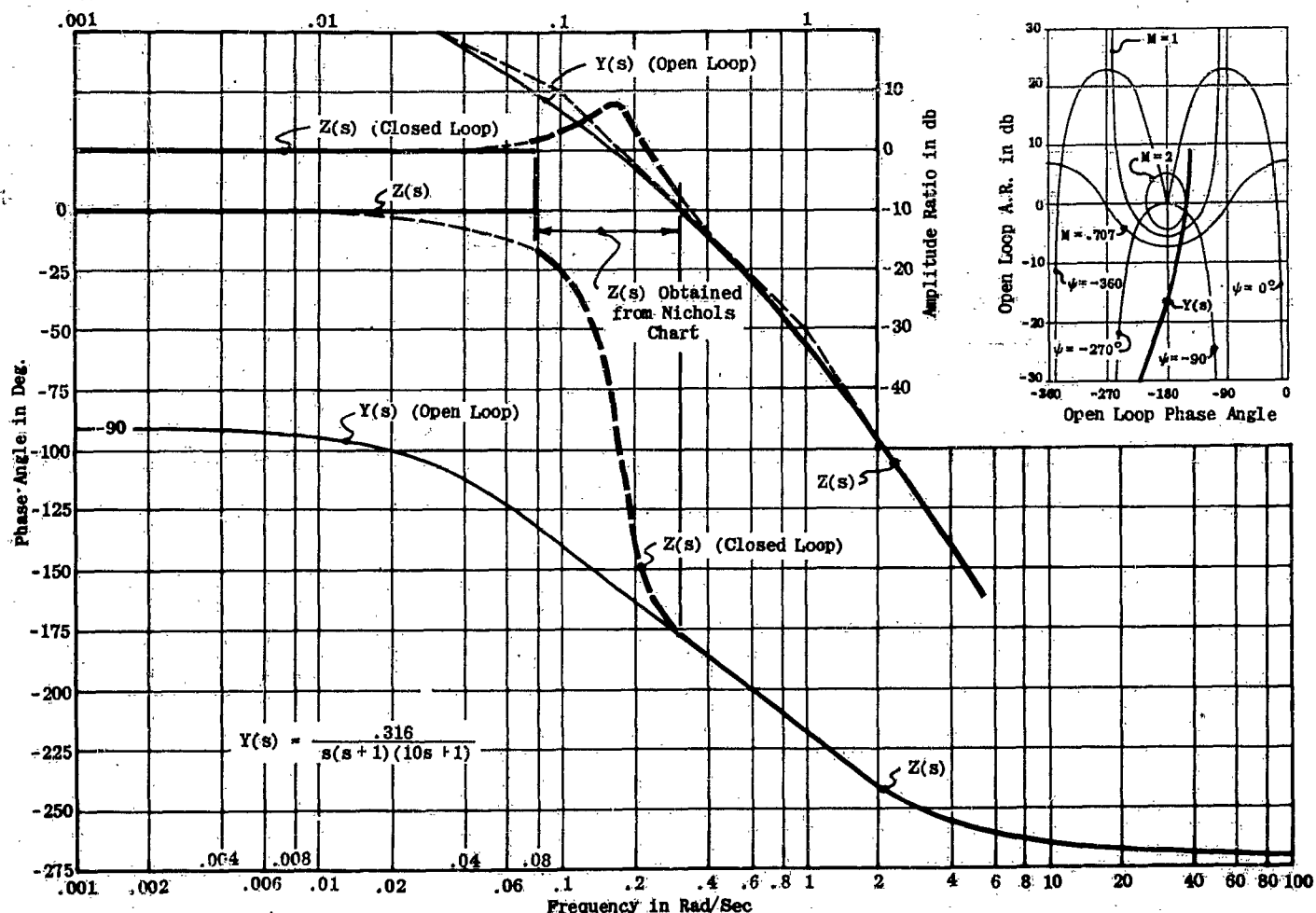


Figure III-17: Open Loop to Closed Loop by Approximations

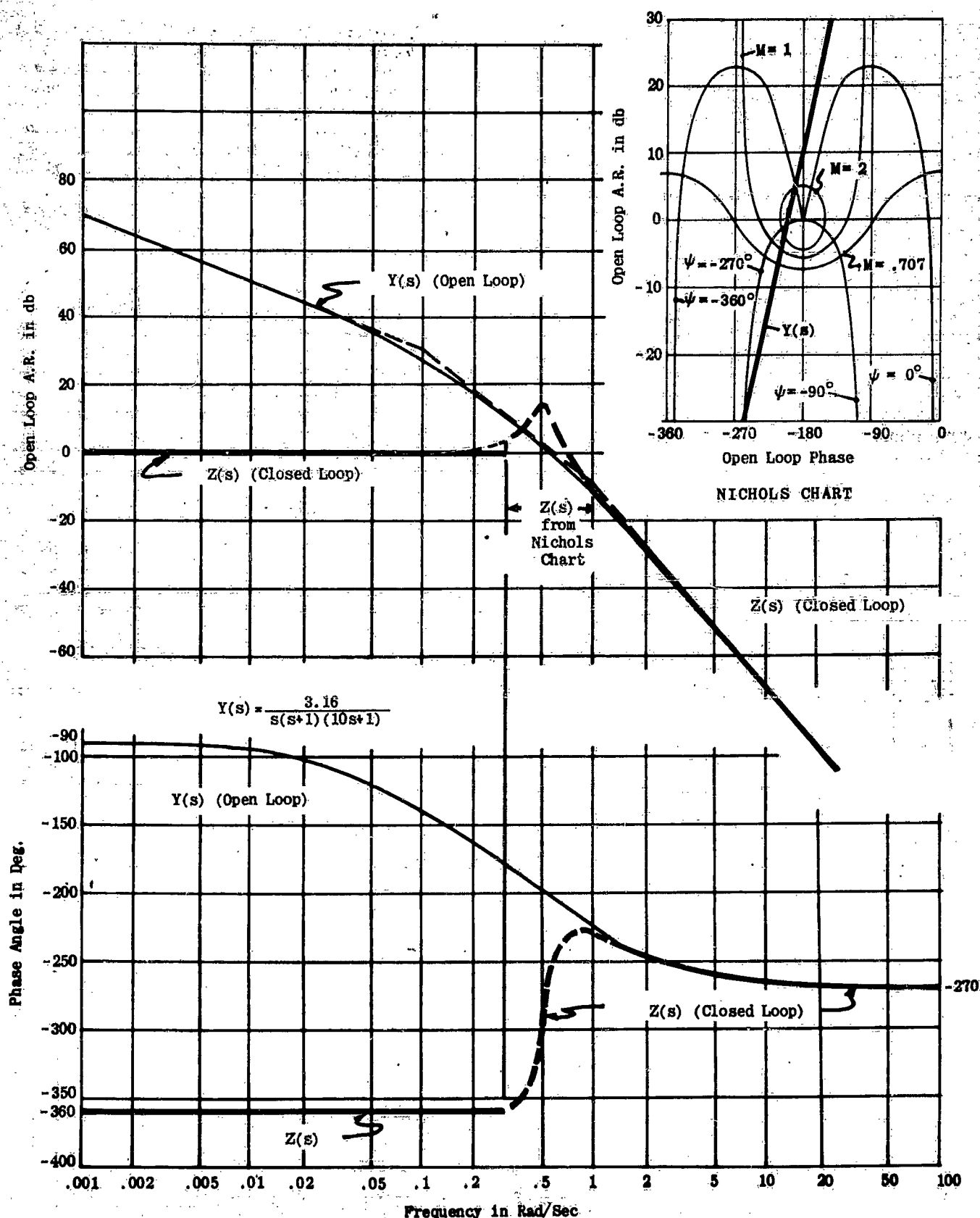


Figure III-18. Open Loop to Closed Loop by Approximations

curve are first drawn. The figure will now be examined by starting from the left (or low frequency side) and continuing to the right (or high frequency side).

- It is evident that at zero frequency the phase is -180° . This indicates that a negative sign (-) should be placed in front of the closed-loop gain K .
- The amplitude decreases at -20 db/dec after the break at $1/\tau_a$. This indicates a first order term in the denominator (or pole). However, the phase lag is decreasing in the vicinity of $1/\tau_a$. This, together with the fact that the first order term is in the denominator, indicates that the term is non-minimum phase of the form $1/(-\tau_a s + 1)$.
- The amplitude breaks from -20 db/dec to 0 db/dec at $1/\tau_b$. This indicates a first order term in the numerator. Since the phase lag continues to de-

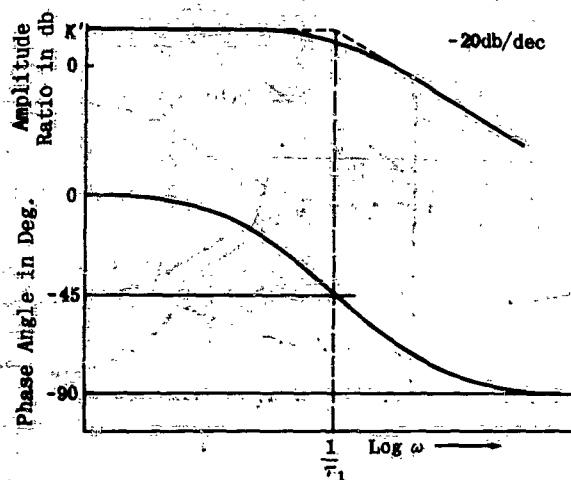


Figure III-19. Obtaining Poles and Zeros of Transfer Function

crease to 0° after this it indicates a minimum phase term of the form $(\tau_b s + 1)$.

- The amplitude breaks to -20 db/dec after $1/\tau_c$ and the phase decreases at the same time, which indicates a first order pole of the form $1/(\tau_c s + 1)$.
- The amplitude breaks from -20 db/dec to -60 db/dec at ω_n . This, together with the peaking at ω_n , indicates a second order term in the denominator. The sharp increase in phase lag at ω_n indicates further that the term is minimum phase. The natural frequency is known from the break point, and the damping may be determined from the amount the peak departs from the asymptotes. (This was discussed in an earlier section). The term then has the form $1/[(s/\omega_n)^2 + 2\zeta(s/\omega_n) + 1]$.

f. The amplitude breaks from -60 db/dec to -40 db/dec at τ_d . This indicates a first order term in the numerator. Since the phase lag increases in the vicinity of $1/\tau_d$ this indicates a non-minimum phase term of the form $(-\tau_d s + 1)$. All the terms are now collected and the closed loop transfer function is of the form

$$Z(s) = \frac{-K(\tau_b s + 1)(-\tau_d s + 1)}{(-\tau_a s + 1)(\tau_c s + 1) \left[\frac{s^2}{\omega_n^2} + \frac{2\zeta s}{\omega_n} + 1 \right]}$$

with the parameters τ , ζ , ω_n obtained from the break points.

In the normal analysis case, the numerator terms τ_b and τ_d are known at the outset, and aid in establish-

ing the asymptotic plot.

EXAMPLE 3. This final example is an exceptional one but is offered to demonstrate that, when the analytic form of the closed-loop transfer function is obtained from the graphical representation, care must be taken

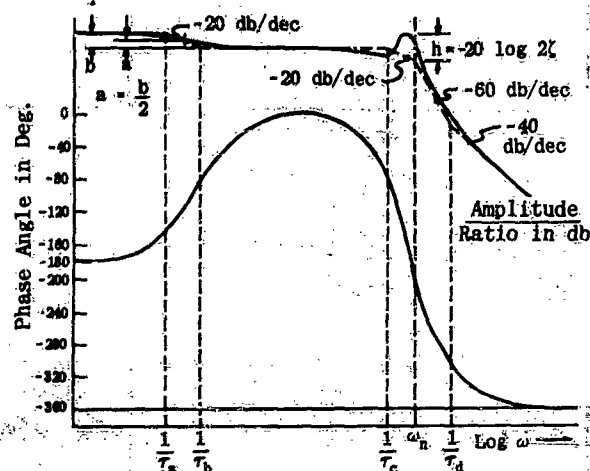


Figure III-20. Obtaining Poles and Zeros of Transfer Function

to consider each and all of the following items:

- The analytic expression for the open-loop transfer function.
- The graphical representation of the amplitude ratio of the closed-loop transfer function.
- The graphical representation of the phase of the closed-loop transfer function.

Consider the closed-loop transfer function shown in figure III-21.

If the amplitude ratio curve alone were considered, the incorrect conclusion might be reached that the transfer function is merely of the form $Z(s) = K = 1$. However, the open-loop transfer function has the form $K(-\tau_a s + 1)/(\tau_b s + 1)$. It is known that the closed-loop transfer function must have the same zeros as the open-loop (i.e., a term in the numerator equal to $(-\tau_a s + 1)$). The next problem is to determine the

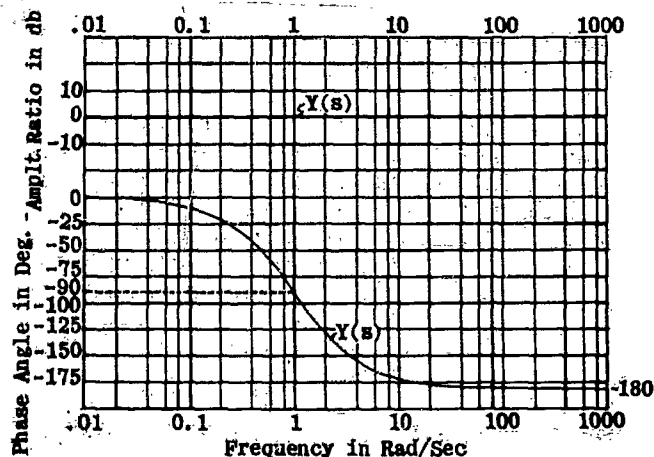


Figure III-21. Obtaining Poles and Zeros of Transfer Function

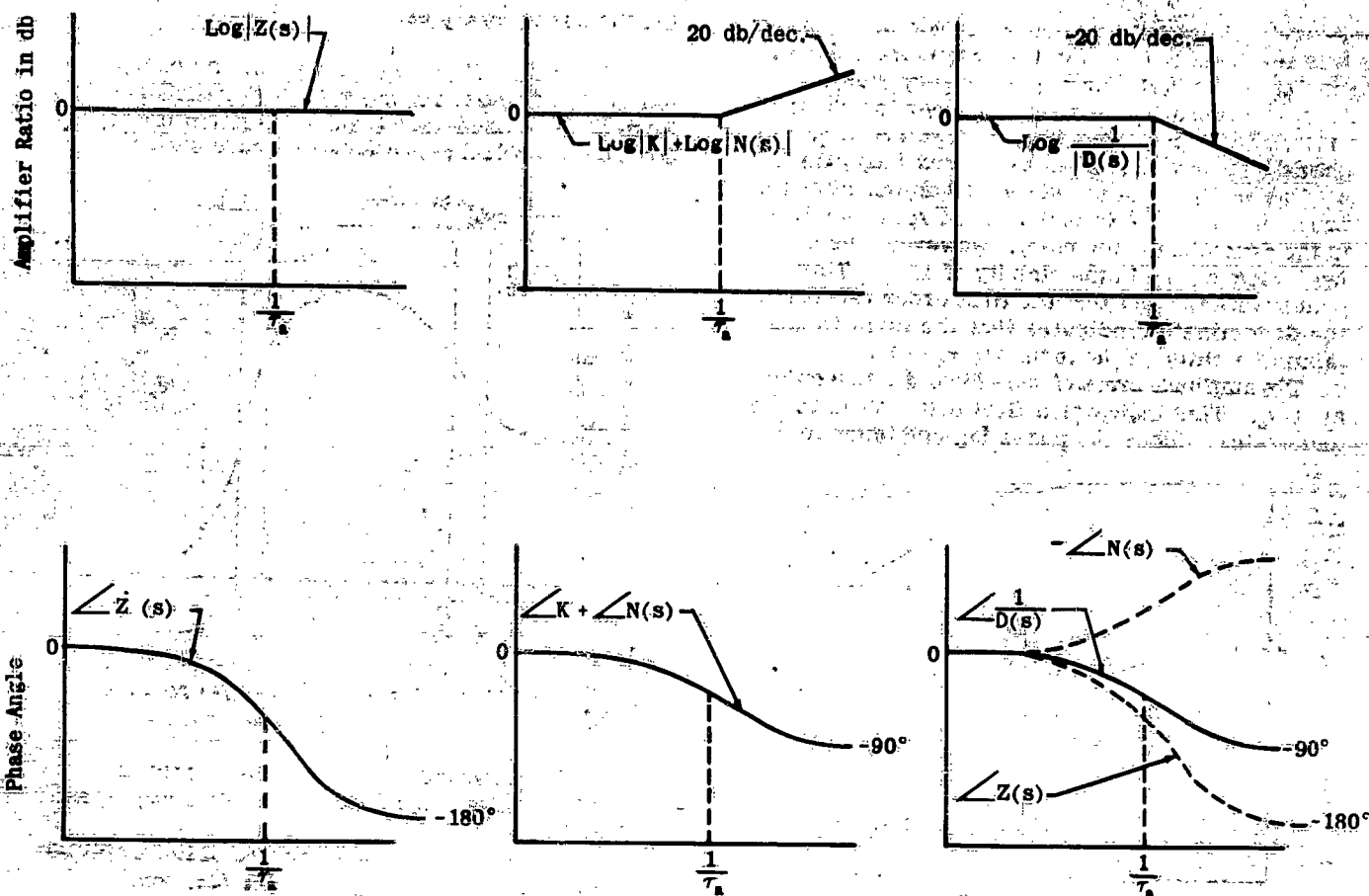


Figure III-22. Graphical Representation of (III-27) and (III-28)

form of the denominator. Since

$$Z(s) = \frac{KN(s)}{D(s)}$$

or $\log|Z(s)| = \log|K| + \log|N(s)| - \log|D(s)|$
(III-27) $\log|D(s)| = \log|K| + \log|N(s)| - \log|Z(s)|$

also

(III-28) $-\angle D(s) = \angle Z(s) - \angle K - \angle N(s)$

When (III-27) and (III-28) are performed graphically as shown in figure III-22, the amplitude and phase of $1/D(s)$ are determined. It is evident that $1/D(s) = 1/(\tau_a s + 1)$ and thus $Z(s) = (-\tau_a s + 1)/(\tau_b s + 1)$. With a little experience the graphical step shown in figure III-22 may be eliminated and the complete analytical form of $Z(s)$ may be written by inspection of its graphical representation.

SECTION 4 — ROOT LOCUS METHOD

(a) INTRODUCTION

Although each of the preceding sections has the purpose of locating the zeros and poles of the closed-loop transfer function, the work was performed not on the s -plane (in which the poles and zeros are located) but on the $Y(s)$ -plane. Of course, this was done as a matter of convenience in obtaining quick approximations to the solution of the problem. The cost of this procedure is that the values of the poles and zeros must be more or less extracted from the graphical representations used.

The root locus method deals with the s -plane exclusively and yields a plot of the locus of all possible roots of $D(s) + KN(s) = 0$ (where $D(s) + KN(s)$ is the denominator of the closed-loop transfer function), as a function of gain. Thus, a mere glance suffices to discover a

nearly complete picture of the dynamics of the system.

(b) BASIC PRINCIPLES

The fundamental problem in establishing a locus of roots of $D(s) + KN(s) = 0$ is the same as existed in the preceding methods, that is, how can these roots be found by working with $Y(s)$ only. To determine the root locus method answer to this question, first note that

(III-29) $1 + Y(s) = 1 + \frac{KN(s)}{D(s)} =$

$$\frac{D(s) + KN(s)}{D(s)} = 0$$

Hence, if $D(s)$ is finite, roots of $D(s) + KN(s) = 0$ are

zeros of $1 + Y(s)$. When $D(s) = \infty$ a special case occurs. In physical problems the order of $D(s)$ is always greater than that of $N(s)$, so that for finite K

$$\lim_{s \rightarrow \infty} \frac{D(s) + KN(s)}{D(s)} = \frac{D(s)}{D(s)} = 1$$

which is not equivalent to (III-29). However, if K is infinitely large

$$\lim_{s \rightarrow \infty} \frac{D(s) + KN(s)}{D(s)} = \frac{\infty}{\infty}$$

and the limiting process might possibly reveal roots.* Consequently, except under this condition, the solution $D(s) = \infty$ is trivial; in all other cases the equation $D(s) + KN(s) = 0$ will reveal all of the zeros of $1 + Y(s)$ of interest (i.e., zeros corresponding to finite values of K). This is the first principle upon which the root locus method is based.

The second important fact upon which the method depends is that the denominator of the closed-loop transfer function of a single loop feedback system is of the form $1 + Y(s)$. Because of this, instead of solving the equation $1 + Y(s) = 0$, it is possible to work with the equivalent expression $Y(s) = -1$. Now, $Y(s)$ can be written as a complex number $R(s)e^{j\phi(s)}$ for any value of s . Also, -1 is a complex number, $-1 = 1e^{j(2k+1)\pi}$. Consequently, if $Y(s) = -1$ then it must be true that:

$$(III-30) \quad Re^{j\phi} = e^{j(2k+1)\pi}$$

where $k = 0, \pm 1, \pm 2, \dots$ so that $k = 1$ and $\phi = (2k+1)\pi$. That is, the magnitude of the complex number, $Y(s)$, must be unity, and its phase angle an odd multiple of π if s itself can be a root of $1 + Y(s) = 0$.

$Y(s)$ is usually of the form:

$$(III-31) \quad Y(s) = \frac{KN(s)}{D(s)}$$

$N(s)$ and $D(s)$ are usually written in the form:

$$s^2(\tau_1 s + 1)(\tau_2 s + 1) \dots \left[\frac{s^2 + 2\zeta_1}{\omega_{n1}^2} s + 1 \right] \left[\frac{s^2 + 2\zeta_2}{\omega_{n2}^2} s + 1 \right] \dots$$

For purposes of working with root loci, it is convenient to rewrite this as:

$$(III-32) \quad \frac{\tau_1 \tau_2 \dots}{\omega_{n1}^2 \omega_{n2}^2 \dots} s^n \left(s \pm \frac{1}{\tau_1} \right) \left(s \pm \frac{1}{\tau_2} \right) \dots$$

$$\times [s^2 \pm 2\zeta_1 \omega_{n1} s + \omega_{n1}^2] [s^2 \pm 2\zeta_2 \omega_{n2} s + \omega_{n2}^2] \dots$$

$$= \frac{\tau_1 \tau_2 \dots}{\omega_{n1}^2 \omega_{n2}^2 \dots} s^n \left(s \pm \frac{1}{\tau_1} \right) \left(s \pm \frac{1}{\tau_2} \right) \dots$$

$$\times [(s \pm \sigma_1 + j\omega_1)(s \pm \sigma_1 - j\omega_1)] [(s \pm \sigma_2 + j\omega_2)(s \pm \sigma_2 - j\omega_2)] \dots$$

where $\sigma_h = -\zeta_h \omega_{nh}$ and $\omega_{nh} = \omega_{nh} \sqrt{1 - \zeta_h^2}$. When $N(s)$ and $D(s)$ are written in this form, $Y(s)$ can be expressed as:

$$(III-33) \quad Y(s) = \frac{KN'(s)}{D'(s)}$$

* This situation is discussed in greater detail later on.

where $N'(s)$ and $D'(s)$ are the factors containing s and the quantity K is the ratio $\tau_1 \tau_2 \dots / \omega_{n1}^2 \omega_{n2}^2 \dots$. Now each of the factors of $N'(s)$ and $D'(s)$ is a complex number (for any value of s) and hence can be plotted as a vector, as in the left hand column of figure III-23.

Since
(III-34)

$$\frac{KN'(s)}{D'(s)} = \frac{[K \kappa r_{N0} r_{N1} \dots r_{N_{n-1}} r_{N_n}] e^{j(\phi_{N0} + \phi_{N1} + \dots + \phi_{N_n})}}{[r_{D0} r_{D1} \dots r_{D_{n-1}} r_{D_n}] e^{j(\phi_{D0} + \phi_{D1} + \dots + \phi_{D_n})}}$$

where the r_{N_i} are the magnitudes and the ϕ_{N_i} the phase angles, of the vectors representing the factors of $N'(s)$; and the r_{D_j} and the ϕ_{D_j} , the corresponding quantities for $D'(s)$.

The quantity in the brackets in (III-34) can be written as:

$$\frac{K \kappa r_{N0} r_{N1} \dots r_{N_{n-1}} r_{N_n}}{r_{D0} r_{D1} \dots r_{D_{n-1}} r_{D_n}} = \frac{K \kappa \prod_{i=0}^n r_{N_i}}{\prod_{j=0}^n r_{D_j}}$$

It is also convenient to denote the total phase angle of $N'(s)/D'(s)$ by

(III-35)

$$\phi_{N0} + \phi_{N1} + \dots + \phi_{N_n} - \phi_{D0} - \phi_{D1} - \dots - \phi_{D_n} = \sum_{i=0}^n \phi_{N_i} - \sum_{j=0}^n \phi_{D_j}$$

Then, by expressing each of the factors of $N'(s)$ and $D'(s)$ as vectors, the conditions that s be a root of $1 + Y(s) = 0$ can be written as:

(III-36)

$$\left| \frac{K \kappa \prod_{i=0}^n r_{N_i}}{\prod_{j=0}^n r_{D_j}} \right| = 1$$

$$\sum_{i=0}^n \phi_{N_i} - \sum_{j=0}^n \phi_{D_j} = (2k+1)\pi \quad k = 0, \pm 1, \pm 2, \dots$$

(when the + sign appears before K)

or

$$\sum_{i=0}^n \phi_{N_i} - \sum_{j=0}^n \phi_{D_j} = 2k\pi \quad k = 0 \pm 1, \pm 2, \dots$$

(when the - sign appears before K)

However, for reasons that will become apparent later, it is more convenient to shift the vectors as in the right hand column of figure III-23. Consequently, roots of $1 + Y(s)$ can be determined by choosing a trial point on the s -plane and drawing vectors to this point from each of the poles and zeros of $Y(s)$.* If the products of the lengths of the vectors r_{N_i} and of $K\kappa$ divided by the product of the lengths of the vectors $r_{D_j} = 1$, (see equation III-34) and the sum of the angles $-(2k+1)\pi$, then the trial point represents a root of the equation $D(s) + KN(s) = 0$ (+ sign before K). Angles measured from a reference line through s are positive

* Note that by shifting the vectors as shown in the right hand column of Fig. (III-23), the point $1/\tau_1$ and $\sigma \pm j\omega$ become the poles or zeros of $Y(s)$. This is one advantage of the shift.

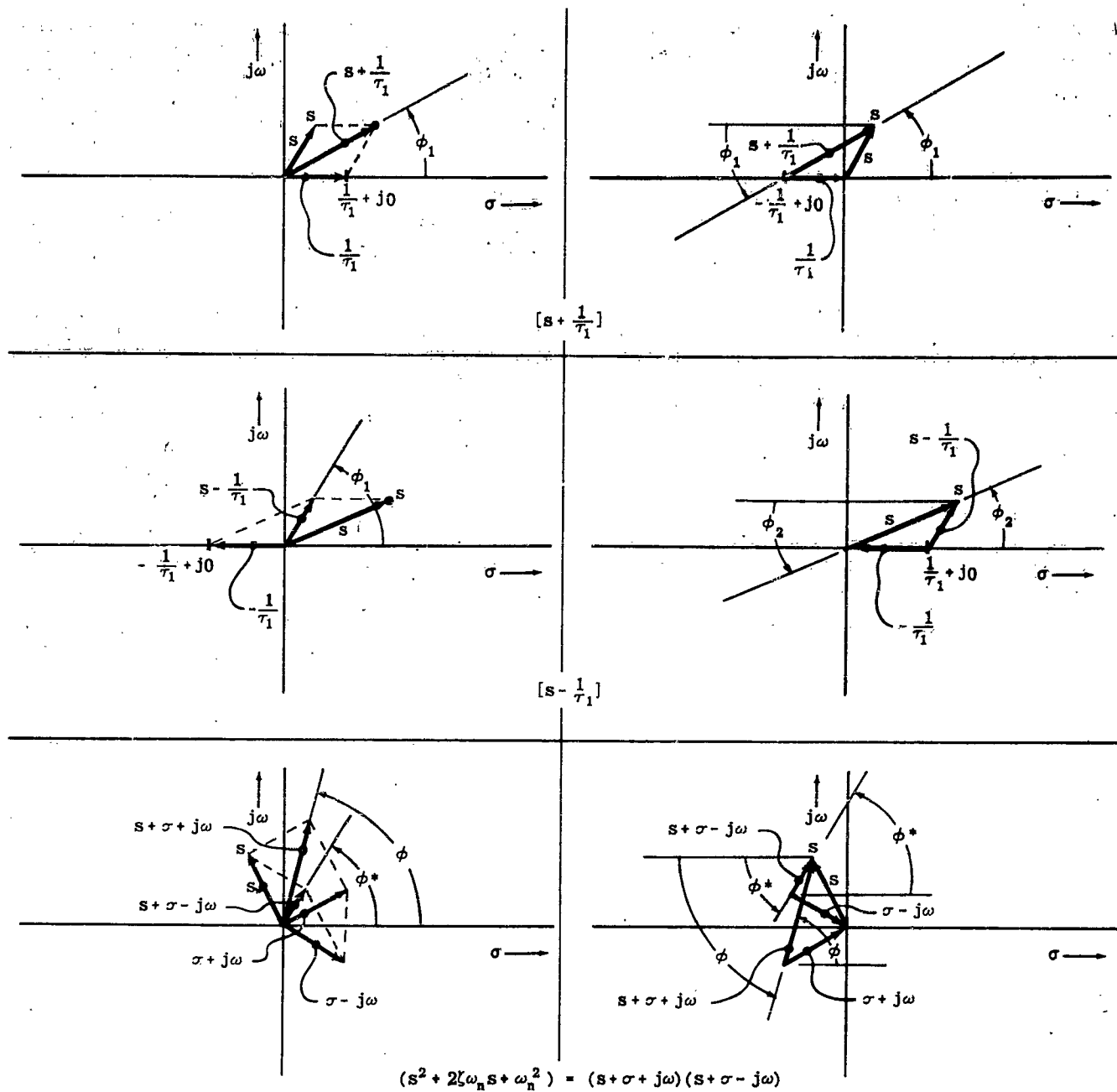


Figure III-23. Phase Angle of Vectors $(s + 1/\tau_1)$, $(s - 1/\tau_1)$, $(s + \sigma + j\omega)$, and $(s + \sigma - j\omega)$

when measured as shown in figure III-24.

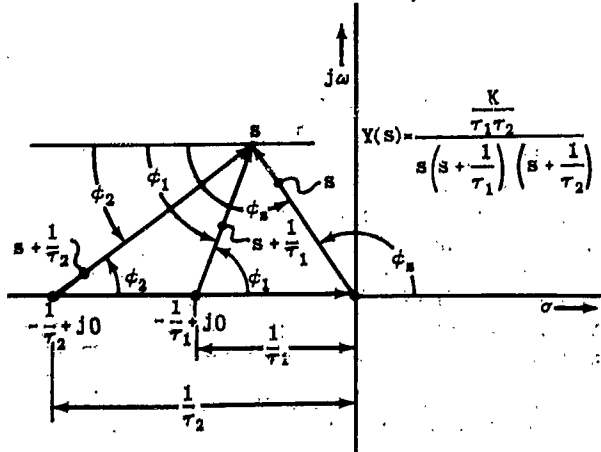


Figure III-24. Typical Construction

It is to be noted that since $\sum \phi \neq (2k+1)\pi$, in this example, the point cannot possibly be a root. Another trial point must be chosen, and the angles added to test its possibilities of being a root. With experience, only very few trials are needed to locate a point that meets the test:

(III-37)

$$\sum_{i=0}^n \phi_{N_i} - \sum_{j=0}^n \phi_{D_j} = (2k+1)\pi \quad k = 0, \pm 1, \pm 2, \pm \dots$$

for + K

or

$$= 2k\pi$$

for - K

The root locus method consists of determining a number of such points in this fashion and drawing a curve through these points, the locus of possible roots; and

for all points along the locus

$$(III-38) \quad K \kappa \frac{\prod_{i=1}^n r_{N_i}}{\prod_{j=1}^m r_{D_j}} = 1$$

must be satisfied. The process is considerably speeded up by the use of a graphical aid (the "Root Locus Plotter") to be described later.

(c) CONSTRUCTING THE LOCUS

As with any graphical construction, root locus plotting is speeded up by establishing:

1. Starting points.
2. End points.
3. Special intermediate points.
4. Asymptotes.

In constructing a locus of roots these points and the asymptotes can be established by inspection.

The starting points of the root locus plot are defined as those points that represent roots of $D(s) + KN(s) = 0$ when K approaches zero. Since the order of $D(s)$ is always greater than that of $N(s)$,

$$(III-39) \quad \lim_{K \rightarrow 0} D(s) + KN(s) = D(s)$$

for any value of s at all. Then the equation $D(s) + KN(s) = 0$ becomes, in this case, $D(s) = 0$. Hence the locus of possible roots of $D(s) + KN(s) = 0$ starts (at gain $K = 0$)* at the zeros of $D(s)$, i.e., the poles of $Y(s)$. As K is increased from zero, the locus must move away from the poles of $Y(s)$ in some direction not yet determined (see figure III-25).

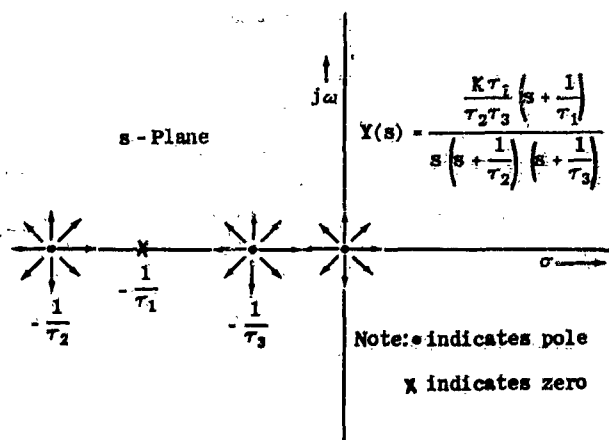


Figure III-25. Starting Points of Locus

The end points are defined as those points that represent roots of $D(s) + KN(s) = 0$ when K approaches infinity. However, it was shown at the beginning of section (b) that there was a possibility of having roots which them-

* When the gain (K) of a closed-loop system is zero an input cannot cause an output. Consequently, the concept of 'roots at zero gain' may be confusing. What actually happens is that as the gain gets very small, the roots approach certain finite values while the amplitudes of the transient terms approach zero. Thus, the type of response is defined by the finite roots of the denominator of $Y(s)$, but the response itself becomes unperceptibly small in magnitude.

selves are very large ($s \rightarrow \infty$) when K is large ($K \rightarrow \infty$). Consequently, there are two conditions which must be investigated to locate the end points:

1. $\lim_{K \rightarrow \infty} [D(s) + KN(s)], \quad (s \text{ finite})$
2. $\lim_{s \rightarrow \infty} [D(s) + KN(s)], \quad (K \text{ large})$

In the first condition

$$(III-40) \quad \lim_{K \rightarrow \infty} D(s) + KN(s) = \lim_{K \rightarrow \infty} KN(s), \quad (s \text{ finite})$$

and, in this event, $D(s) + KN(s) = 0$ reduces to $KN(s) = 0$. Consequently, at high gain (K) the locus must enter each of the zeros of $N(s)$, figure III-26, which are also the zeros of $Y(s)$.

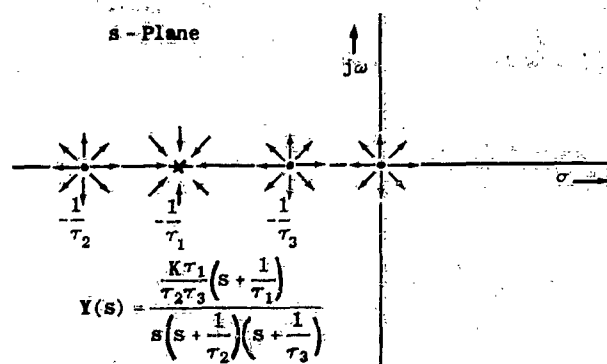


Figure III-26. Starting and End Points of Locus

But since the order of $N(s)$ is always less than that of $D(s)$, and since the order of $D(s) + KN(s)$ is equal to the order of $D(s)$, there are $(n-m)$ * roots still to be accounted for at high gain. Since condition (1) did not reveal them, these remaining roots must correspond to condition (2). That is, the roots occur at $s \rightarrow \infty$. This is illustrated in figure III-27. An observer located far away from the origin of the s -plane and from all the zeros and poles sees the zeros and poles of $Y(s)$ all bunched up on the s -plane near the origin.

In fact, from far enough out in the plane, the vector angles of all the factors become very nearly equal. Then the angle of a vector from a zero to the trial point far out on the plane appears to cancel the angle of a vector from one of the poles to the trial point. The result is that each zero appears to cancel a pole so that the plot of $Y(s)$ looks like simply a multiple order pole at the origin. The order of the pole will be the difference between the number of poles and zeros of $Y(s)$. That is

$$(III-41) \quad \lim_{s \rightarrow \infty} Y(s) = \lim_{s \rightarrow \infty} \frac{KN(s)}{D(s)} = \frac{K}{s^{n-m}}$$

Consequently for large values of s there appear to be no zeros of $Y(s)$, and $D(s) = s^{n-m}$. So that

$$(III-42) \quad \lim_{s \rightarrow \infty} D(s) + KN(s) = s^{n-m}$$

* n and m are the orders of $D(s)$ and $N(s)$ respectively. See equation (III-36).

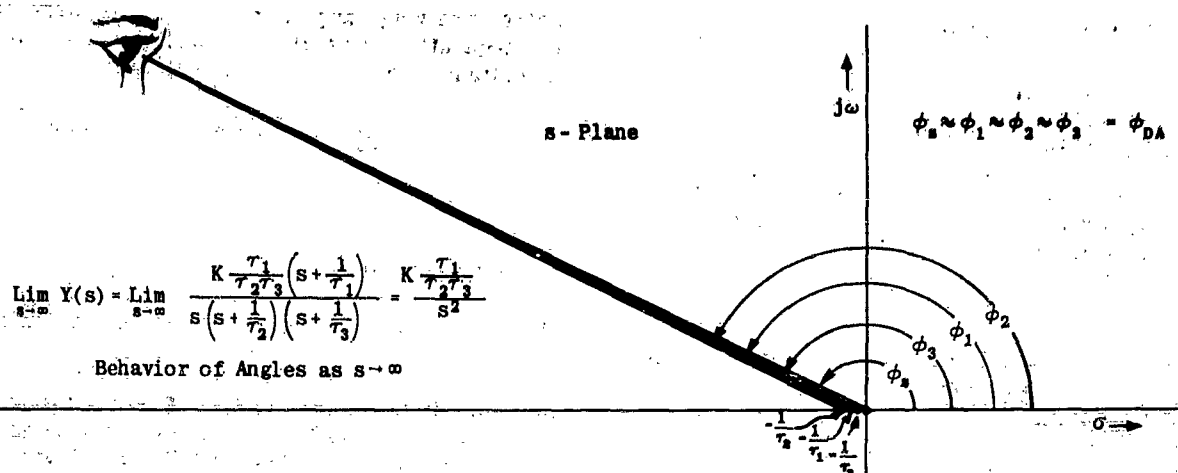


Figure III-27. Limiting Values of Angles as $|s| \rightarrow \infty$

Now, since $\sum_{i=1}^n \phi_{N_i} = 0$ and $\phi_{D_1} = \phi_{D_2} = \dots = \phi_{D_j} = \dots = \phi_{D_A}$

$$\sum_{i=1}^n \phi_{N_i} - \sum_{j=1}^m \phi_{D_j} = -\sum_{j=1}^m \phi_{D_A} = -(n-m)\phi_{D_A} = (2k+1)\pi \text{ or } 2k\pi$$

$$k = 0, \pm 1, \pm 2, \pm 3, \dots$$

then

$$(III-43) \quad \phi_{D_A} = \frac{(2k+1)\pi}{n-m} \text{ or } \frac{2k\pi}{n-m} \quad k = 0, \pm 1, \pm 2, \pm 3, \dots$$

Evidently vectors drawn pointing away from the origin at angles ϕ_{D_A} will point to large roots of $D(s) + KN(s) = 0$ for large values of K .

At this point it is well to recall that it has been determined that at zero gain the roots of $D(s) + KN(s) = 0$ are simply the roots of $D(s) = 0$ and that as gain (K) increases, the locus must proceed away from these roots. It has also been determined that under some conditions the roots of $D(s) + KN(s) = 0$ for large values of K are the roots of $KN(s) = 0$ and for other cases they appear far out in the s -plane at angles ϕ_{D_A} . Since the locus has several starting points, and several terminating points, it would appear that the locus is not a single continuous curve, but rather consists of several branches. Certain branches will start at the poles of $Y(s)$ (roots of $D(s) = 0$) and proceed to zeros of $Y(s)$ (roots of $KN(s) = 0$) and other branches will start at other poles and proceed toward the roots at infinity in the direction of the angles ϕ_{D_A} . Thus it would appear that (III-43) establishes asymptotes for $n-m$ branches of the locus. It will be shown in the following pages that this is true (see figure III-28).

Starting points, end points, and asymptotes have been established. Some intermediate points can be established by closer examination of certain parts of the s plane. Consider first the real axis. To determine whether or not there are any possible roots there, take a trial point (s) on the real axis, figure III-29.

All the vectors from the zeros and poles to point s_1 make angles nearly zero. Consequently $\sum \phi \approx 0$ and there cannot possibly be any roots along the positive real axis. Next try a point s_1 near the real axis between $-1/\tau_3$ and the origin (figure III-30). In this case the vector from the pole at the origin has an

Since there is a + sign before K , angles to asymptotes are

$$\phi_{DA_1} = \frac{180^\circ}{2} = 90^\circ \quad \phi_{DA_2} = \frac{(3)(180^\circ)}{2} = 270^\circ \quad n-m=2$$

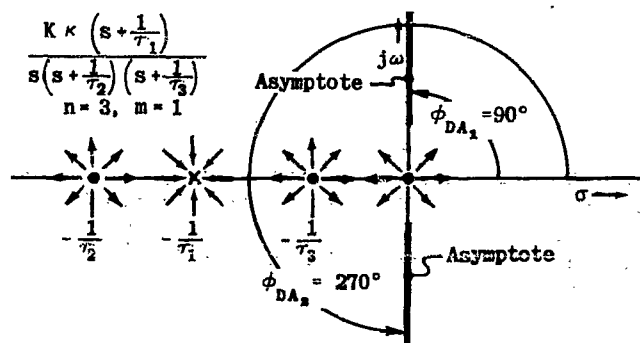


Figure III-28. Starting Points, End Points and Asymptotes

angle of approximately 180° and the rest of the poles and zeros contribute nothing. Consequently, there can be roots anywhere between the pole $-1/\tau_3$ and the pole at the origin. It should be noted that the

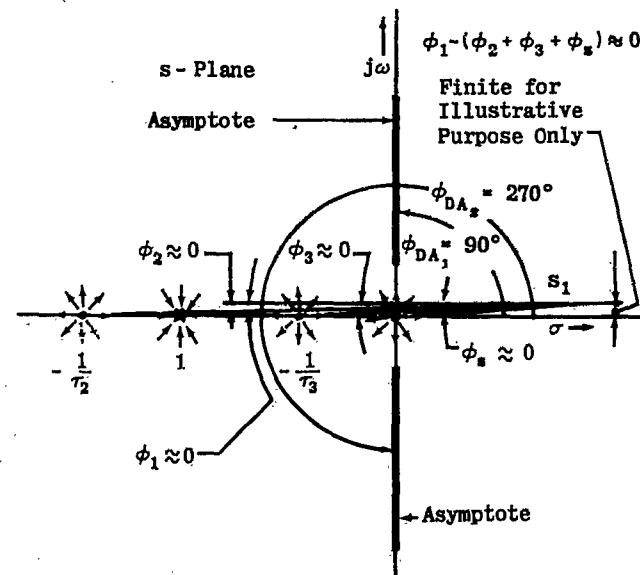


Figure III-29. Test for Roots on Positive Real Axis

Chapter III Section 4

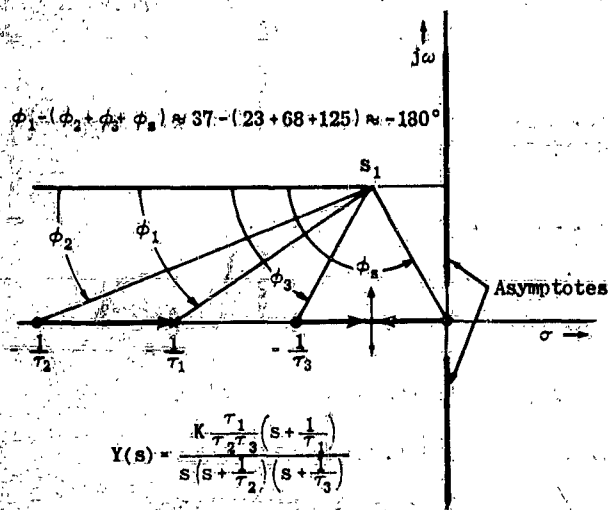


Figure III-33. Trial Angle Calculations

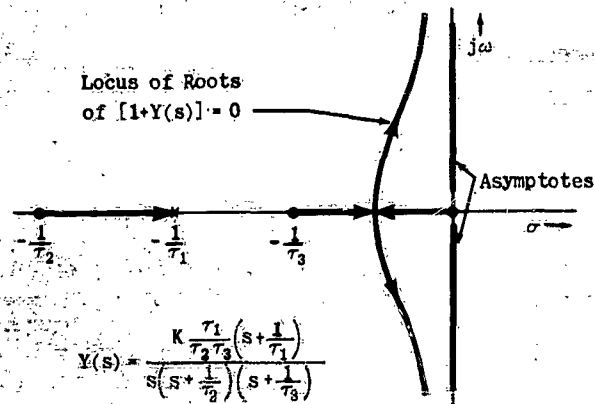


Figure III-34. Complete Locus

In order to locate the roots corresponding to a known value of gain (K), points on each of the branches of the locus are chosen and the lengths of the vectors measured. If these lengths satisfy (III-36), the points represent roots for this value of K . However, in the more common problem K is unknown and is to be determined in order to place the roots at certain desired locations as discussed in section II-3e. In this problem, the desired root on the locus is chosen and the gain (K) computed as the unknown in equation (III-36).

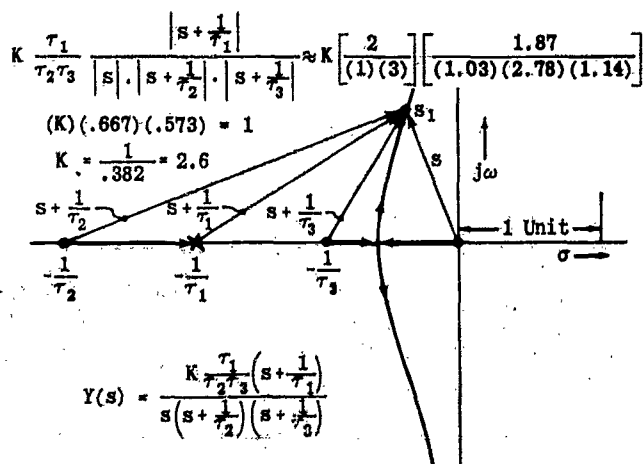


Figure III-35. Calculation of K for a chosen Root s_1

When $Y(s)$ contains complex roots, the procedure is identical to that previously described. However, since

the graphical construction has a slightly different appearance the following discussion is included to avoid possible confusion. The example also includes a pole on the positive real axis.

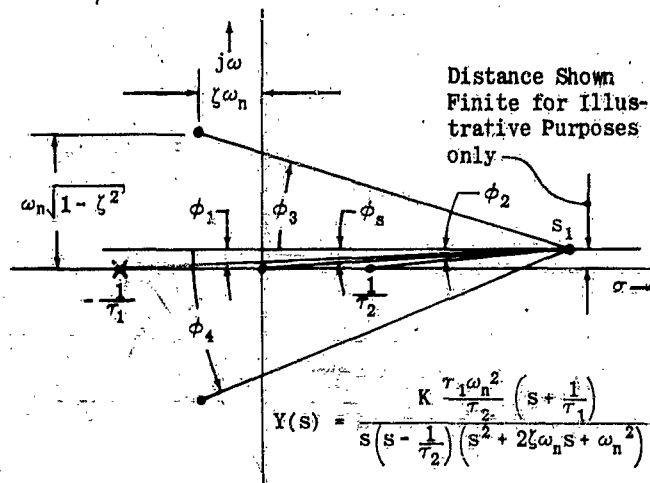


Figure III-36. Location of Locus on Real Axis

In figure III-36 a point near the positive real axis beyond $+1/\tau_2$ is chosen to check for roots there. The angles of the vectors from the poles and zeros on the real axis are all zero. The angles to vectors originating at the complex poles are non-zero but they are equal and opposite so the net angle is zero; therefore there are no roots beyond the $1/\tau_2$. Repeating this procedure along the real axis locates roots between the origin and the pole at $1/\tau_2$ and between $-1/\tau_1$ and $-\infty$. Consequently, the locus along the real axis appears as in figure III-37. This illustration also shows

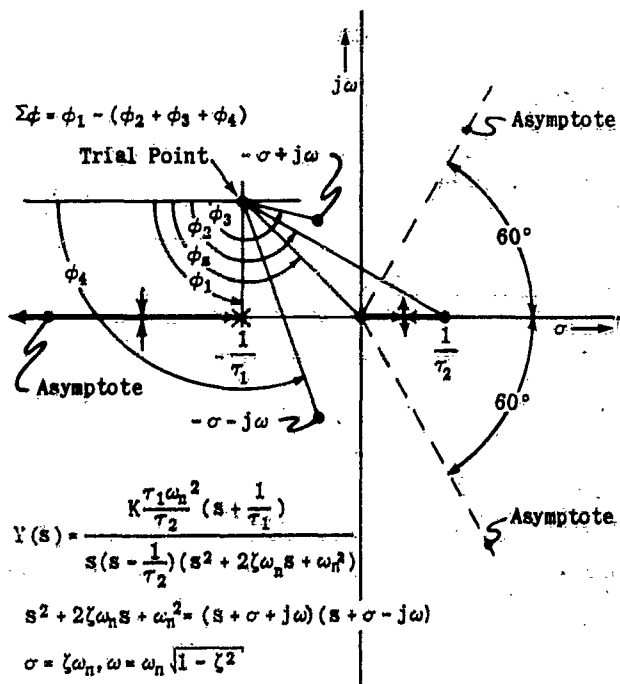


Figure III-37. Construction Involving Complex Poles

the asymptotes and how the vectors are drawn from the complex poles. (The points at which the locus enters and leaves the real axis are still undetermined but are shown to indicate what the locus must do.)

In order to determine in which direction the locus moves away from the complex poles, a trial point is taken very close to one of them. The point is taken so close to the pole that the vectors drawn from all the other poles and zeros to the trial point look as though they terminated right on the complex pole being investigated. This is shown in figure III-38. Although this trial point is taken very near the complex pole its exact location above, below, left or right of

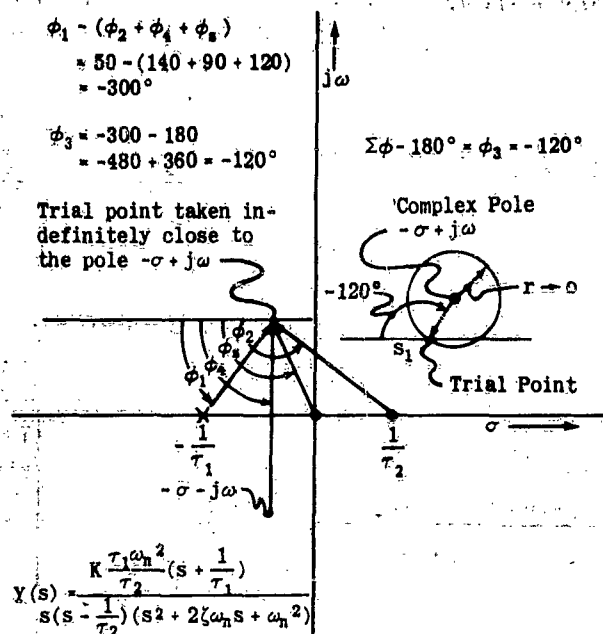


Figure III-38. Establishing Direction of Departure from Complex Pole

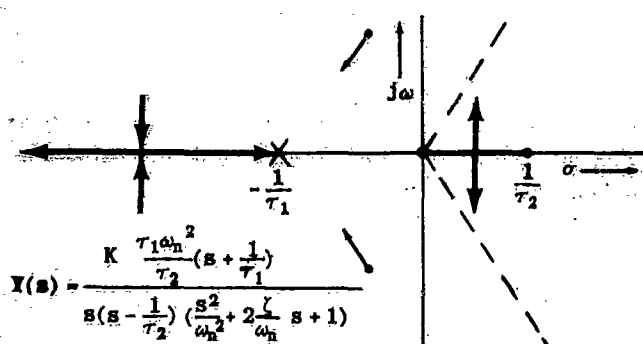


Figure III-39. Initial Parts of Locus

the pole is left undetermined. Consequently, when the angles are added, the angle to the vector from the complex pole being investigated to the trial point is left out (in this example ϕ_3); and in order that the trial point lie on the locus, the sum of the angles in this example must be $\phi_1 - (\phi_2 + \phi_4 + \phi_5) = 180K_{\text{odd}}$ or $\phi_3 = \phi_1 - (\phi_2 + \phi_4 + \phi_5 + 180K_{\text{odd}})$. In the figure it is seen that $\phi_3 = -120^\circ$. This means that the vector from the complex pole to the trial point lies down and to the left as shown by the inset in figure III-38. With this added information, figure III-39 can be constructed.

Figure III-40 indicates some possible forms of the locus. Several intermediate points must be determined in order to show which is the correct locus.

(d) THE ROOT LOCUS PLOTTER

Constructing a locus of roots using the techniques described would be a tedious job and the method would not be very valuable to the designer. However, it is possible to devise a simple tool that greatly simplifies and expedites the work. The use of this tool will now be discussed.

One form of the plotter is shown in figures III-41 and III-42.

The technique of using the plotter is indicated in figures III-43 through III-52. The pintle is first placed on the trial point. In this position, the card and disc are held together by the action of the spring. The disc and pintle are then depressed by pressing on the top of the pintle. This releases the card from the disc. The card is now free to rotate without carrying the disc with it, the latter being held in its initial position by the pressure on it.

For the purpose of showing clearly what is happening, the plotter will now be reduced to schematic form, figure III-44.

The system shown in figure III-24 will be used for illustration. To determine if a point s_1 lies on the locus of roots for this system, place the pintle of the plotter s_1 as shown in figure III-45.

The disc is depressed and the reference line on the card moved until it coincides with the line from s_1 to $-1/\tau_2$ as shown in figure III-46.

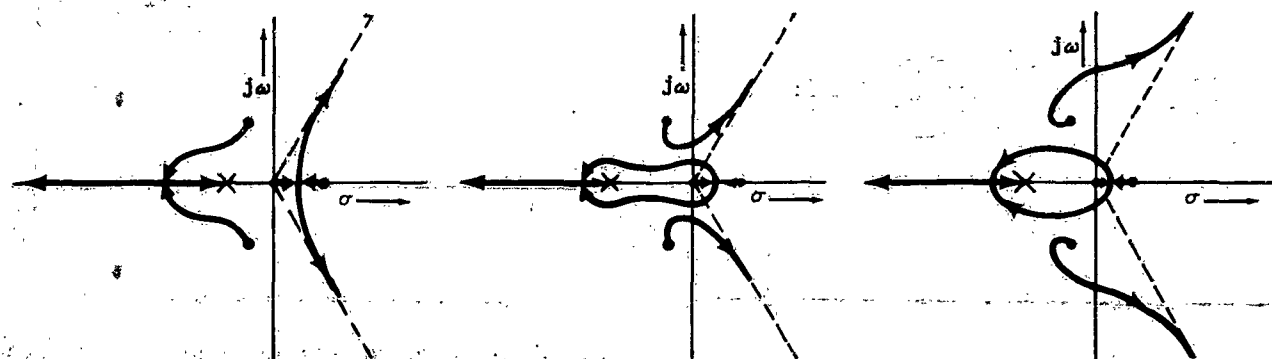


Figure III-40. Root Locus Forms

Chapter III
Section 4

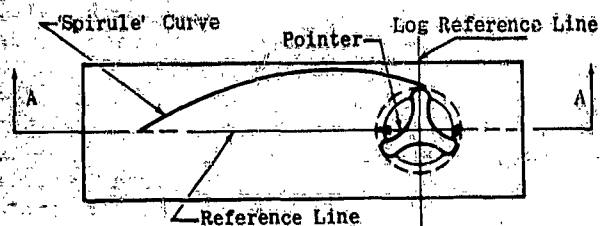


Figure III-41. Diagrammatic Representation of Root Locus Plotter

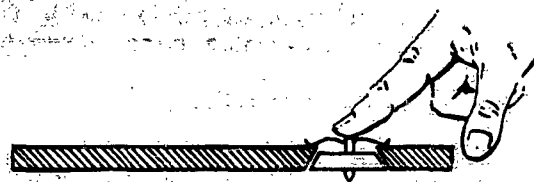


Figure III-43. Cross Section of Root Locus Plotter Showing Card Freed From Disc

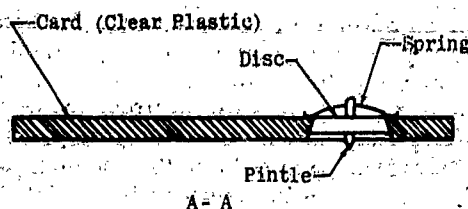


Figure III-42. Cross Section of Root Locus Plotter Showing Disc Locked to Card

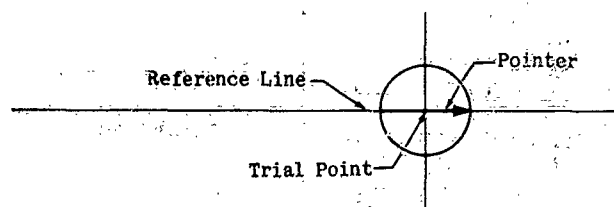


Figure III-44. Schematic of Root Locus Plotter

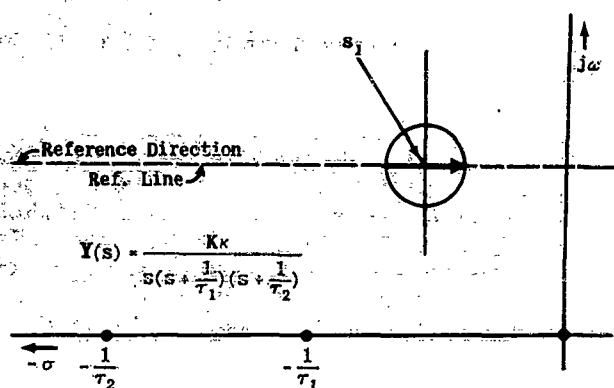


Figure III-45. Use of Root Locus Plotter

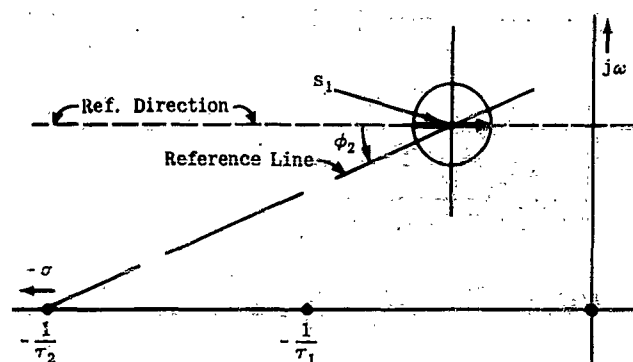


Figure III-46. Use of Root Locus Plotter

The pointer is now released and the entire plotter is swung back until the reference line again coincides with the reference direction, i. e., the horizontal direction through s_1 . This operation will have rotated the disc through an angle ϕ_2 , with respect to the card. The plotter now appears as shown in figure III-47.

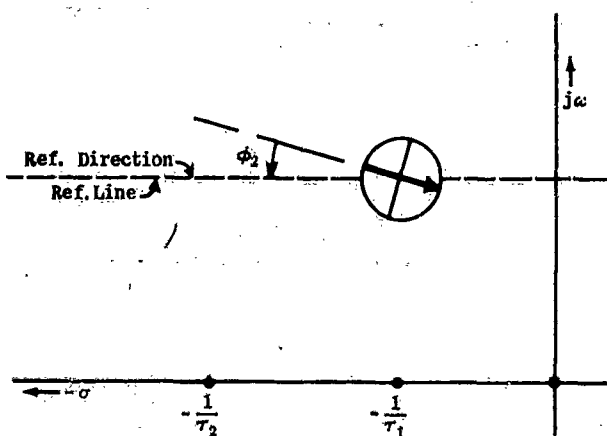


Figure III-47. Use of Root Locus Plotter

This process is now repeated using $-1/\tau_1$, as shown in figures III-48 and III-49.

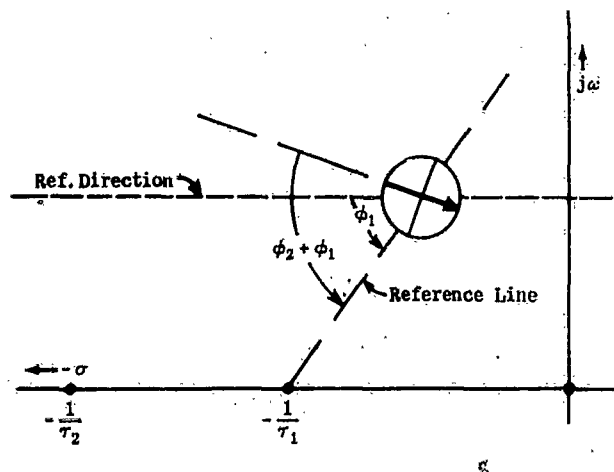


Figure III-48. Use of Root Locus Plotter

The plotter is again reoriented upon the reference line as shown in figure III-49 and now has $\phi_2 + \phi_1$ added into it as shown by the relative position of disc to card.

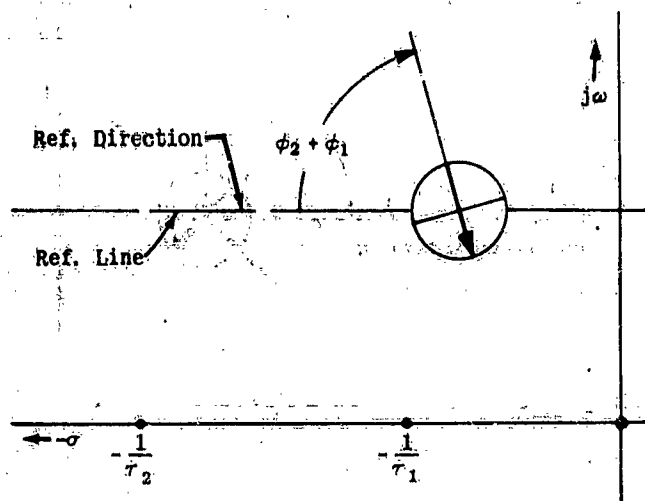


Figure III-49. Use of Root Locus Plotter

The process is repeated for the root at the origin as shown in figures III-50 and III-51.

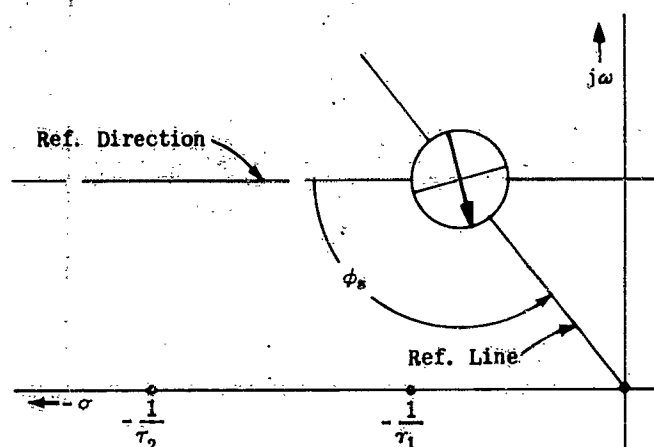


Figure III-50. Use of Root Locus Plotter

When the plotter is returned to its initial position, as shown in figure III-51, it is seen that the sum of the angles thus measured is not equal to 180° . This is indicated by the fact that the pointer has not exactly reversed its direction.

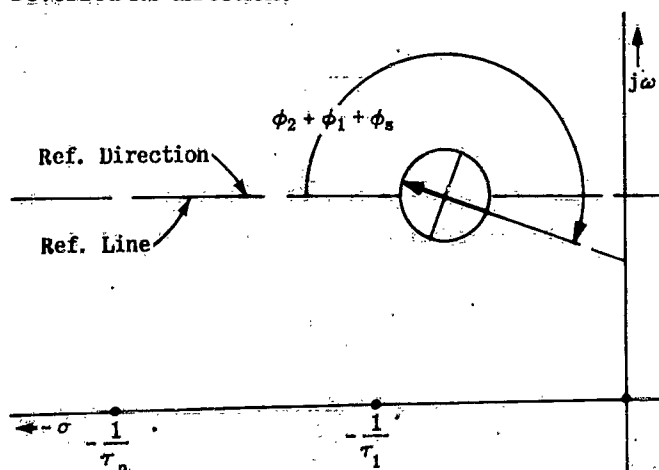


Figure III-51. Use of Root Locus Plotter

Since for this particular $s_1, \phi > 180^\circ$, a second trial point, s_2 , lying on the same horizontal line is selected and the entire process repeated for it. A possible result is shown in figure III-52.

For $s_2, \phi < 180^\circ$. Since the two points, s_1 and s_2 , lie on the same horizontal line, the locus must pass between them.

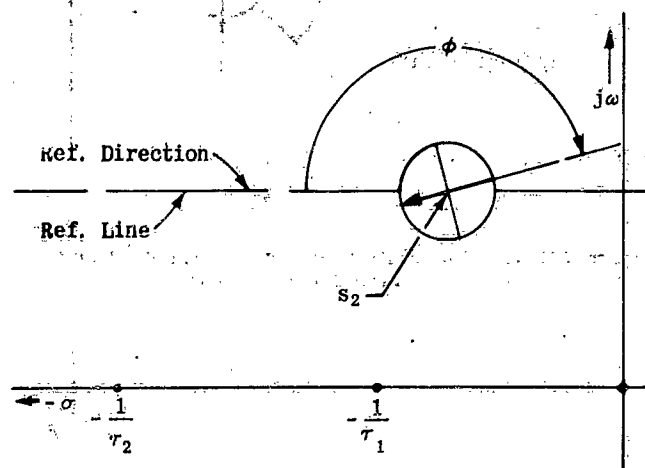


Figure III-52. Use of Root Locus Plotter

When zeros occur in the transfer function it is necessary to subtract their angles from those of the poles. This is done by reversing the order of operations in using the plotter, as is illustrated in the following example, using the transfer function of (III-43).

$$(III-43) \quad Y(s) = Kk \frac{s + \frac{1}{\tau_1}}{s(s + \frac{1}{\tau_2})}$$

The poles and zeros of (III-43) are plotted in figure III-53.

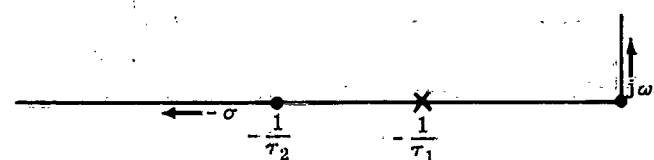


Figure III-53. Poles and Zeros of Equation (III-43)

Select a point s_1 , and add into the plotter the angle to the pole $-1/\tau_2$, ϕ_2 , as shown in figure III-54.

The plotter is returned to its original position as shown in figure III-55.

The next critical point to be considered is $-1/\tau_1$. Since it is a zero, its angle must be subtracted from the sum of the angles of the poles. This is done as follows: Without depressing the pointer, the card is oriented with the reference line lying along the vector joining s_1 and $-1/\tau_1$ as shown in figure III-56.

The pointer is then depressed and the reference line swung through the clockwise angle ϕ_1 to the reference direction as shown in figure III-57.

Chapter III Section 4

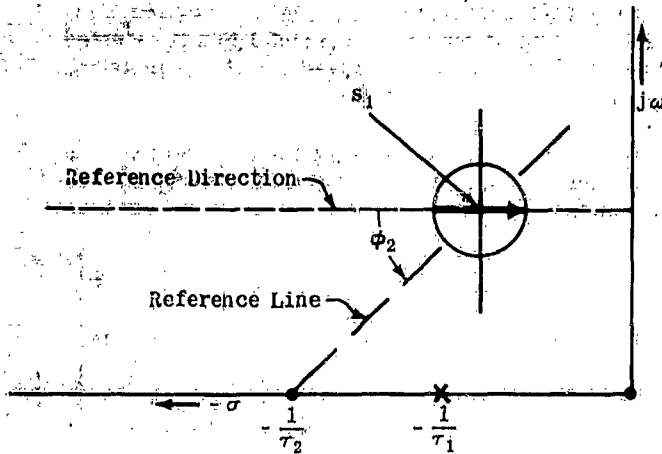


Figure III-54. Use of Plotter for Open-Loop Transfer Function with Zeros and Poles

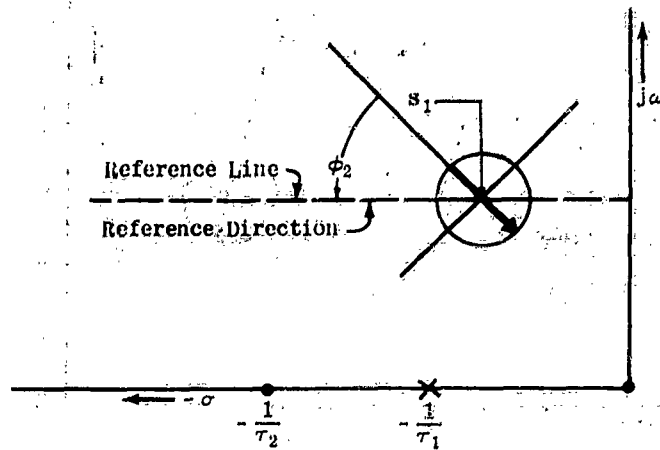


Figure III-55. Use of Plotter for Open-Loop Transfer Function with Zeros and Poles

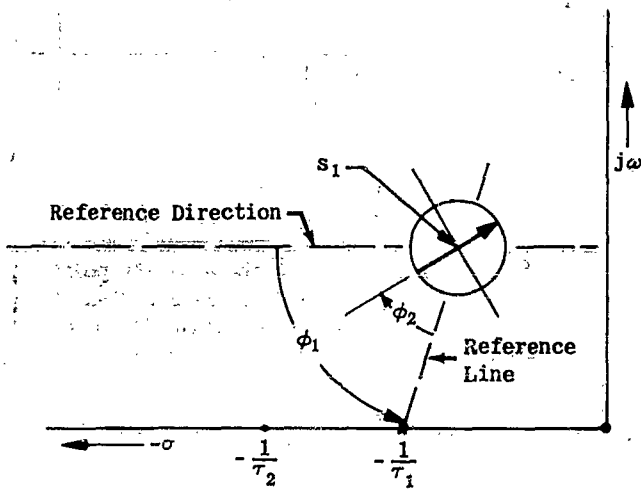


Figure III-56. Use of Plotter for Open-Loop Transfer Function with Zeros and Poles

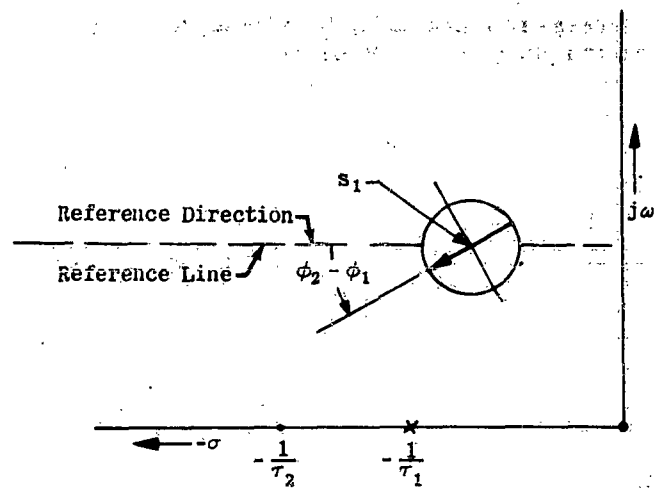


Figure III-57. Use of Plotter for Open-Loop Transfer Function with Zeros and Poles

The angle ϕ_2 that was in the plotter prior to this last step has been reduced by ϕ_1 .

In this way, by rotating the plotter before or after releasing the disc, it is possible to take care of both zeros and poles.

A logarithmic spiral marked on the plotter can be used to multiply the lengths of the vectors according to equation (III-36). The equation of a logarithmic spiral is of the form

$$(III-46) \quad r = a^{\theta/\theta_0}$$

(see figure III-58) or

$$(III-47) \quad \log_{10} r = (\log_{10} a) \frac{\theta}{\theta_0}$$

In these relations, the quantities a and θ_0 are parameters which can be used to set the linear and angular scales.

Consider, as an example, the vectors r_1 and r_2 , such that the corresponding angles on the logarithmic spiral

are θ_1 and θ_2 . Then if $r = r_1 r_2$ (see figure III-59),

$$(III-45) \quad \log_{10} r = \log_{10} r_1 + \log_{10} r_2 = \frac{\log_{10} a}{\theta_0} (\theta_1 + \theta_2)$$

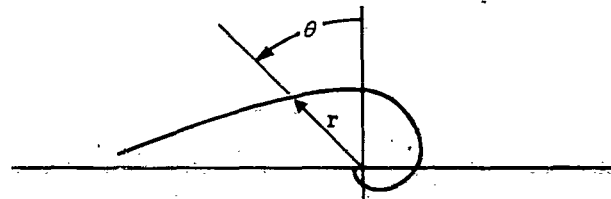


Figure III-58. Logarithmic Spiral

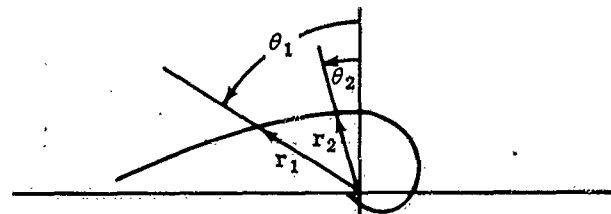


Figure III-59. Multiplication

It is convenient to take $a = 10$, and $\theta_0 = \pi/2$. When this is done, the radius vector of the curve is increased by a factor 10, if the angle θ is increased by 90° .

Figure III-60 shows the general appearance of the part of the logarithmic spiral lying in the second, third, and fourth quadrants. (It is to be noted here that the angle θ is measured from the upward vertical direction, and not from the horizontal.)

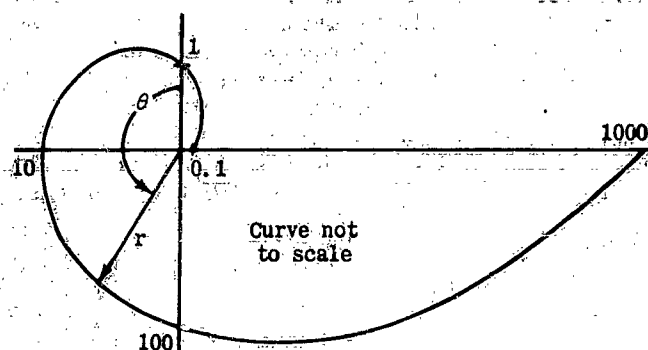


Figure III-60. Multiplication Factors

Figure III-60 shows that if the vector being measured is in the second quadrant, its numerical value lies between 1 and 10; in the third, between 10 and 100; in the fourth, between 100 and 1000, etc. Thus, if the radius vector representing a product of the lengths falls in the second quadrant, that is, if its end lies on the portion of the curve actually inscribed on the plotter, the product is read off directly.

It is a property of this logarithmic spiral that if a radius vector, r_2 , leads another radius vector, r_1 , by 90° , then $r_2 = 10r_1$; for a lead angle of 180° , $r_2 = 100r_1$; and, in general, if r_2 leads r_1 by $n \times 90^\circ$, $r_2 = 10^n r_1$. This makes it possible to find the magnitude of any vector by using only the part of the spiral lying in the second quadrant; if r_2 falls in the third quadrant, read off the value of the radius vector in the second quadrant which lags r_2 by 90° and multiply this value by 10; if r_2 falls in the fourth quadrant, read the value of r in the second quadrant which lags behind r_2 by

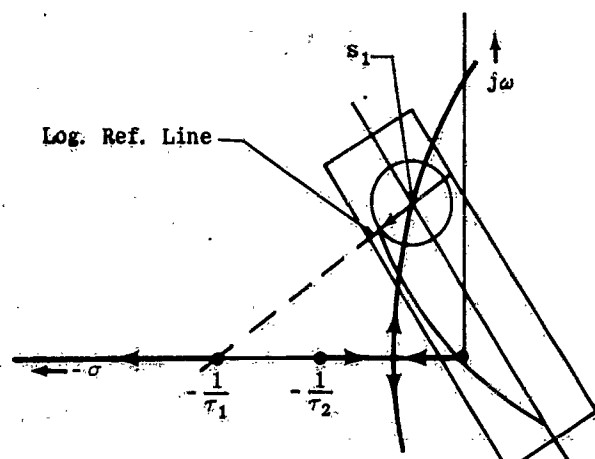


Figure III-61. Determination of K for a Point on the Root Locus of $Y(s) = (K\kappa)/[s(s + \frac{1}{\tau_1})(s + \frac{1}{\tau_2})]$

$2 \times 90^\circ$ (this is simply the radius vector lying along the extension of r_2 backwards, into the second quadrant) and multiply its value by 100, and so on. r is always evaluated in the second quadrant, and the result multiplied by a power of ten whose exponent is the smallest number of clockwise right angles needed to bring the actual radius vector into the second quadrant.

For example, assume that the gain is desired at the point s_1 on the locus of figure III-61. The plotter is to be aligned initially as shown in figure III-61 with the log reference line along the line joining s_1 and $-1/\tau_1$.

The disc is depressed and the card rotated until the log curve falls upon the point $-1/\tau_2$ as shown in figure III-62.

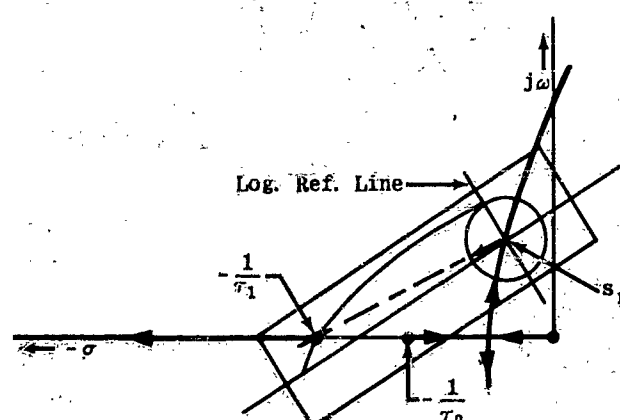


Figure III-62. Determination of K for a Point on the Root Locus of $Y(s) = (K\kappa)/[s(s + \frac{1}{\tau_1})(s + \frac{1}{\tau_2})]$

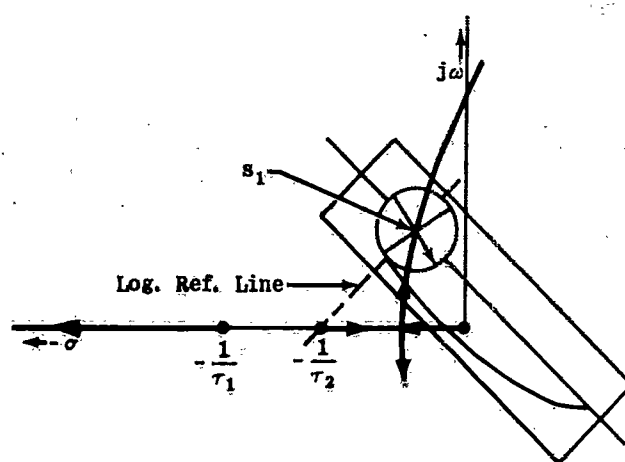


Figure III-63. Determination of K for a point on the Root Locus of $Y(s) = (K\kappa)/[s(s + \frac{1}{\tau_1})(s + \frac{1}{\tau_2})]$

Since the head of the arrow lies in the second quadrant, $|s_1 + 1/\tau_1|$ is measured directly on the logarithmic spiral.

To multiply $|s_1 + 1/\tau_2|$ by $|s_1 + 1/\tau_1|$ the entire plotter is then aligned, with the "log reference line" lying along the vector from s_1 to the root $-1/\tau_2$ as shown in figure III-63.

Chapter III

Section 4

Again the disc is depressed and the card is swung until the log curve falls upon the root $-1/\tau_2$ as shown in figure III-64.

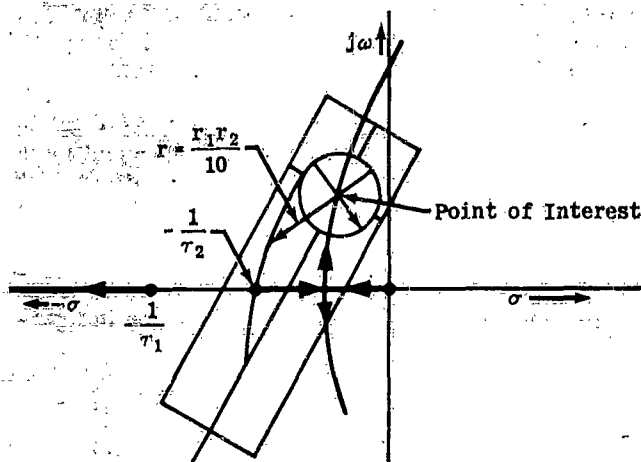


Figure III-64. Determination of K for a Point on the Root Locus

For the roots chosen for this example, the arrow on the disc now lies in the third quadrant; so the product $r_1 r_2$ is read from the plotter as shown and the answer multiplied by 10.

If there are zeros in a $Y(s)$ being investigated, it is still possible to determine the gain in one continuous process; the zeros and poles may be taken into account in any convenient order. It is only necessary to subtract the angle on the spirule corresponding to a zero. This may be done in the same way as the angles due to zeros were subtracted in determining the root locus itself.

(e) APPLICATIONS

The root locus method yields a plot of the poles of the closed-loop transfer function $[Y(s)]/[1+Y(s)]$ as K varies from 0 to ∞ . Consequently, if the locus moves into the right half plane, the possibility of instability exists. In design problems, the locus is usually plotted for a proposed system and the gain is adjusted so that a certain damping ratio (ζ) or

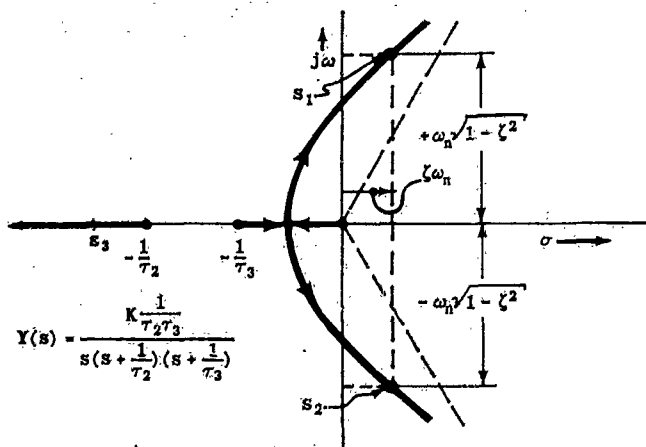


Figure III-65. Root Locus of Simple Systems

damping factor ($\zeta\omega_n$) is obtained. If desired stability cannot be achieved for any gain, equalization becomes necessary (i.e., adding poles and zeros to alter the loci).

Whenever a zero is added to $Y(s)$, it has the effect of drawing the locus toward it. Conversely, whenever a zero is removed, the locus moves away from the vacated point. For example, the system of figure III-34 can never become unstable because of the zero at $-1/\tau_1$. If the zero is removed, the locus appears as in figure III-65. At high values of gain this system becomes unstable because of the roots $s_{1,2} = \zeta\omega_n \pm j\omega_n\sqrt{1-\zeta^2}$ and $s_1 = \zeta\omega_n - \omega_n\sqrt{1-\zeta^2}$ in the right half plane. This cannot happen as long as the zero is present in the left half plane. This is illustrated by figure III-66 in which the zero is located at two different positions.

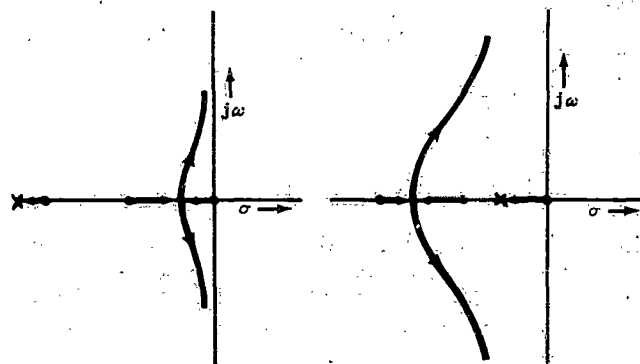


Figure III-66. Effect of Zero on Negative Real Axis

When a complex conjugate pair of zeros is added to the system of figure III-65, the locus appears as in figure III-67, while the zero on the positive real axis draws the locus to itself producing a system that is unstable for any value of gain (see figure III-68).

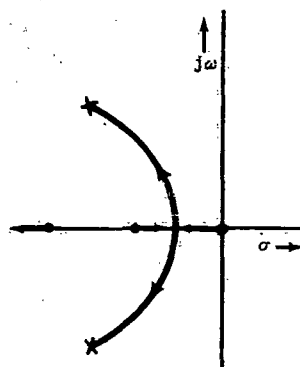


Figure III-67. Effect of a Pair of Complex Conjugate Zeros

Poles repel the locus; this is shown by a resketch of figure III-65 in which the pole at $-1/\tau_2$ is moved along the real axis from the origin to $-\infty$ (see figure III-69).

Figure III-70 illustrates a common occurrence in the analysis of complex systems. As the zero is moved along the real axis a situation is reached in which the branches of the locus have a pair of complex conjugate roots in common. It is not possible to determine this point directly by the methods just discussed, but it

may be obtained by extrapolating the parts of the locus in the immediate vicinity.

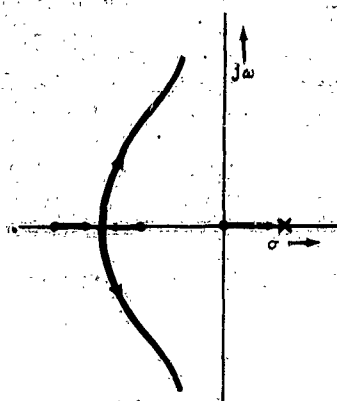


Figure III-68. Effect of Zero on Positive Real Axis

However, if for some reason it is important to know the value of the equal roots (common point in figure III-70(c)), the method described in reference 5 (page 159) can be worked out on the root locus plot. However, this is very rarely of interest.

Section II-3(e) discusses a graphical method of obtaining the coefficients of an equation of motion in the time domain from a plot of the poles and zeros of the equation in the s plane. Now, in the closed-loop transfer function, $Z(s) = [Y(s)]/[1 + Y(s)] = [KN(s)]/[D(s) + KN(s)]$, the zeros are the roots of $KN(s) = 0$ and are known, and the poles are the roots of $D(s) + KN(s) = 0$ which can be obtained on the s -plane by the root locus method. Consequently, the graphical procedure of section II-3 (e), can be carried out right on the root locus plot to obtain the coefficients of the closed-loop equation of motion in the time domain.

A final note of caution should be sounded here concerning the form of $Y(s)$ to be used in the root locus method.

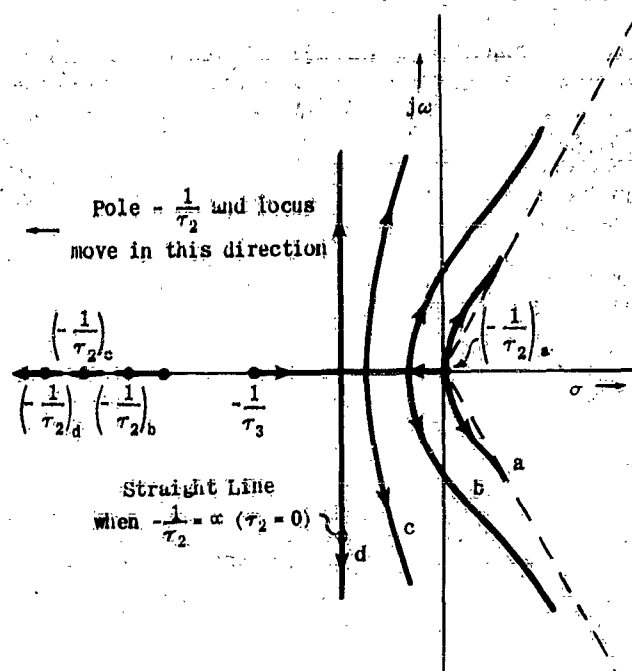


Figure III-69. Effect of Pole Location on Shape of Locus

It is important that in the expression for $Y(s)$

$$Y(s) = \pm K \kappa \frac{(s + \frac{1}{\tau_1})(s + \frac{1}{\tau_2})(s^2 + 2\zeta_1\omega_{n_1}s + \omega_{n_1}^2)}{(s + \frac{1}{\tau_3})(s + \frac{1}{\tau_4})(s^2 + 2\zeta_2\omega_{n_2}s + \omega_{n_2}^2)}$$

all of the s that stand alone (the s of $s + 1/\tau_1$, the $s^2 + 2\zeta_2\omega_{n_2}s + \omega_{n_2}^2$) be preceded by a $+$ sign and that K be a positive number. Unless this is done consistently, the analyst will frequently use the criterion $\Sigma\phi = (2k+1)\pi$ when he should be using $\Sigma\phi = 2k\pi$, and worthless results will be obtained.

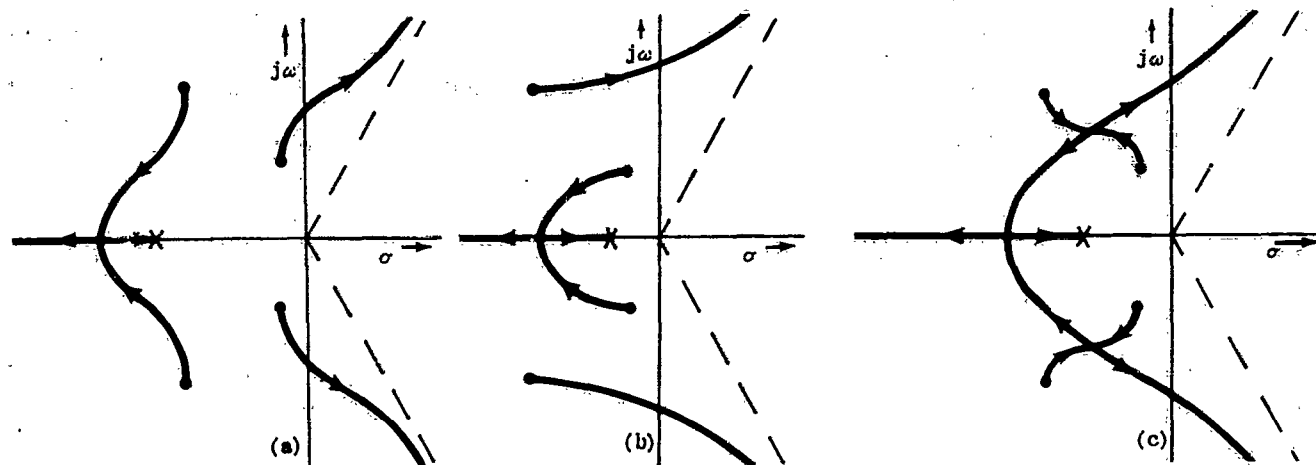


Figure III-70. Effect of Zero Location on Shape of Locus

SECTION 5—SPECIAL CASES—GAIN MARGIN, PHASE MARGIN, MAXIMUM MAGNIFICATION RATIO

While the previous portions of this chapter have covered the most important analytical tools used in servo-

mechanisms work, completeness requires the consideration of certain other concepts that are often

Chapter III Section

useful. These concepts are all concerned with obtaining approximate data regarding the closed loop transfer function, and hence the transient response. The approach used in this section shall be to define pertinent quantities initially and then to discuss what information may be obtained from a knowledge of these quantities.

An open loop transfer function polar plot is shown in figure III-71.

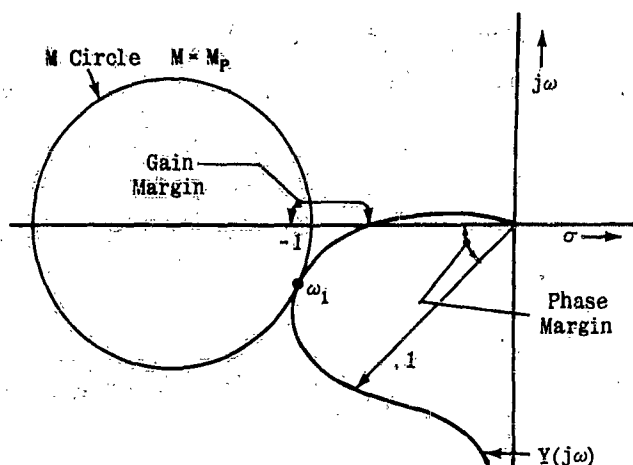
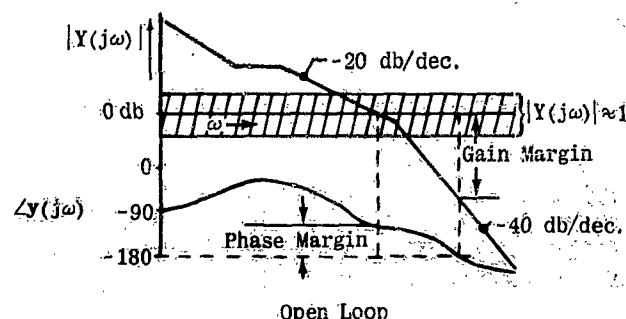


Figure III-71. M_p , and Gain and Phase Margins

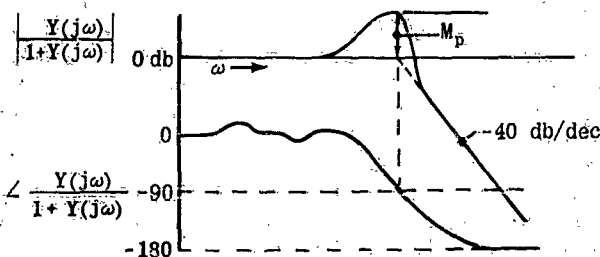
The first of the quantities to be considered is the maximum value of the closed loop transfer function, which is defined by the M circle just tangent to the open loop locus. This value of M is called the peak magnification ratio of the closed loop system, and is denoted by M_p . The second quantity is the phase margin; it is the angle between the negative real axis and the $Y(j\omega)$ vector for the frequency at which $|Y(j\omega)|$ is unity. The third quantity is the gain margin, and is the value of $|1 - Y(j\omega)|$ for the frequency at which $Y(j\omega)$ has a phase angle of -180° . (The same quantities could of course have been defined in the same way utilizing an open-loop logarithmic plot. The polar-plot was used only for convenience). It should be noted that each of the quantities approximately defines the closeness of the open-loop plot to the minus one point. Since the close-

ness of the open-loop transfer function to this point is a measure of relative stability of the modes existing in this region, the M_p value, phase margin, and gain margin have been used extensively as measures of stability (hence the term "margin"). In systems of certain specific types, criteria in terms of these quantities have served as measurements of both stability and transient response. However, their use in such a fashion for any system in general without extensive investigation is likely to yield misleading results.

In this subsection, a slightly different point of view will be taken.



Open Loop



Closed Loop

Figure III-72. Approximation by Second Order System

The phase margin, gain margin, and M_p value serve as bounds to the open-loop transfer function in the region of frequencies giving the largest values of the closed-loop transfer function, (i.e., where $|Y(j\omega)| \approx 1$). It is to be expected, therefore, that they can be correlated with the poles and zeros of the closed-loop

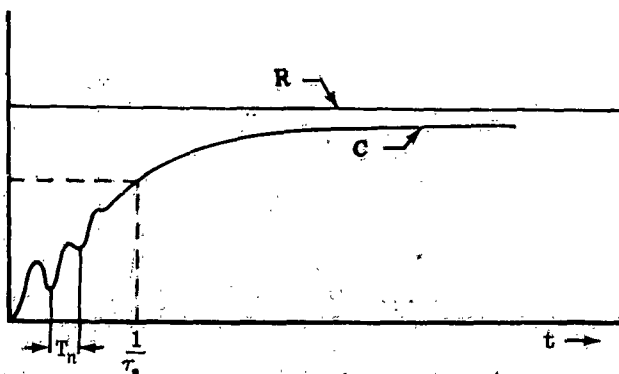
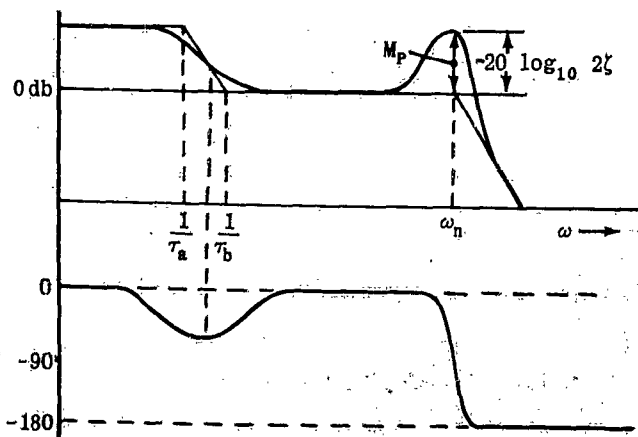


Figure III-73. Special Transfer Function

transfer function occurring near this region. If such a correlation can be determined, then the region where $|Y(j\omega)|$ is of the order of 1 can be approximated along with those regions defined by (III-24) and (III-25) (i.e., $|Y(j\omega)| \ll 1$ and $|Y(j\omega)| \gg 1$). The closed-loop transfer function can then be reasonably well known for all values of ω with a mere glance at the open loop function.



Figure III-74. Commercial Model Root Locus Plotter

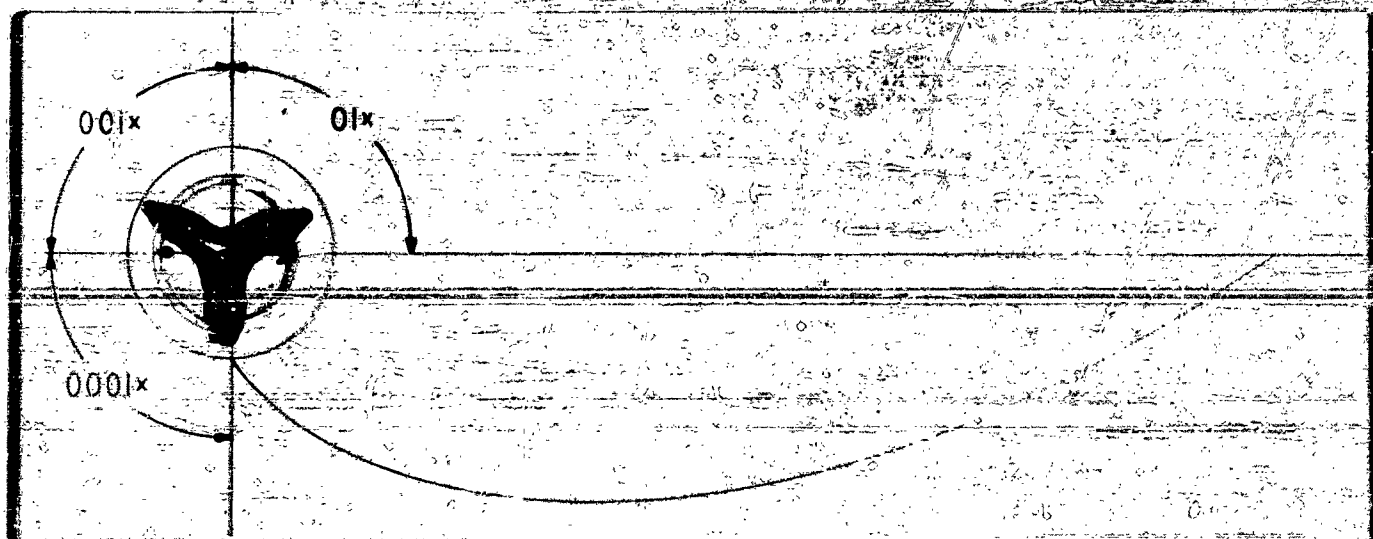


Figure III-75. Northrop Model Root Locus Plotter

If the servo system is such that only one break from a -20 db/dec slope to a -40 db/dec slope occurs in the region where $|Y(j\omega)| \approx 1$, the closed loop transfer function in that region is normally approximated by a quadratic factor (see figure III-72).

In this case, the closed loop transfer function looks like that of a second order system. In a second order system, the maximum magnification ratio is given by:

$$(III-49) \quad M_p = \frac{1}{2\zeta\sqrt{1-\zeta^2}}$$

and occurs at a frequency

$$(III-50) \quad \omega = \omega_n \sqrt{1-2\zeta^2}$$

where ζ = damping ratio and ω_n = undamped natural frequency. (III-27) is plotted (in db) in figure A-13, and (III-28) is plotted in figure A-12. Examining the

transient responses of second order systems, figures A-1 and A-2, it is noted that the overshoot tends to become excessive for values of $\zeta < 0.4$. For good transient response, therefore, the value of M_p for such a system should be from 1.2 to 1.6. (1.7 db to 4 db).

Since many servos are reasonably close to that of figure III-72 in the regions where $|Y(j\omega)| \approx 1$, it is a widespread practice to use the second order system peak (M_p) as a specification of transient response. The phase margin can be discussed in a similar fashion, with representative values of 30 to 40° being considered reasonable. For gain margin, values of 0.6 to 0.9 (-4.5 db to -1 db) are often quoted as being adequate.

There are many exceptions to the accurate use of M_p , gain margin, or phase margin criteria in predicting transient response of systems. For example, a closed loop transfer function of the type shown in figure III-73 has the form (if minimum phase).

The M_p value in this case can be used only to estimate the response time of the quadratic mode, whereas the actual response time will be determined almost entirely by the value of τ_a if M_p is a reasonable value.

$$(III-51) \quad Z(s) = \frac{K(\tau_b s + 1)}{(\tau_a s + 1)\left(\frac{s^2}{\omega_n^2} + \frac{2\zeta}{\omega_n}s + 1\right)}$$

It is logical to conclude however, that M_p , gain margin, and phase margin criteria give reasonable answers about the transient response in the immediate region

of the modes they partially define. If (III-24) and (III-25) are also used, a fair idea of transient response can then be obtained in many practical cases with comparatively little effort. However, since the effort required to make a much more complete analysis by the open loop-closed loop, or even the root locus method, is but little more, it is recommended that these more exact methods be used for the first solutions of any new system. From this first solution, M_p , gain margin and phase margin criteria can be defined and used to minimize the effort required for the following analyses.

BIBLIOGRAPHY

The following bibliography is included for reference. The list is in no sense complete, but contains the major source material for this chapter. Many of the references, themselves, contain much more complete and detailed bibliographies.

1. 'Theory of Servomechanisms,' by James, Nichols, and Phillips; McGraw Hill Book Co., 1947.
2. 'Network Analysis and Feedback Amplifier Design,' by H. W. Bode; D. Van Nostrand Co., New York, 1945
3. 'Servomechanisms and Regulating System Design,' by H. Chestnut and R. Mayer; John Wiley and Sons, New York, 1951.
4. 'Principles of Servomechanisms,' by G. S. Brown and D. P. Campbell; John Wiley and Sons, New York, 1948.
5. 'Control System Synthesis by Root Locus Method,' by W. R. Evans; Trans. AIEE, Vol. 69, 1950.
6. 'A Generalization of Nyquist's Stability Criteria,' by A. Vazsonyi, Journal of Applied Physics, Vol. 20, 1949.
7. 'Analysis of Feedback Systems by the Open Loop-Closed Loop Logarithmic Method,' by D. T. McRuer and R. Zacharias, Northrop Aircraft, Inc., Unpublished paper, 1952.
8. 'Introduction to Higher Algebra,' by Maxime Bocher; Macmillan, 1907.

CHAPTER IV SYNTHESIS

SECTION 1 - INTRODUCTION

(a) GENERAL

As a general term, synthesis may be defined as the process of combining elements into a unified whole. A servo system may be considered as the combination of two basic portions — a controlled element and a controller. The controlled element is characterized by output quantities to be controlled and input quantities to which control is applied. The controller has three functions, namely, sensing, actuation, and equalization. The first of these is performed by sensors, or elements capable of detecting the output quantities to be controlled. The second function requires actuators, or elements capable of applying control. The third function, equalization, includes all of the means required to connect or modify the performance of any of the system elements and of the overall system to achieve satisfactory system operation. In summary, the basic functional portions of a servo system are:

1. The controlled elements
2. The controller elements
 - a. Sensing elements (sensors)
 - b. Actuating elements (actuators)
 - c. Equalization

The most general servomechanism system synthesis problem is then one of designing properly both controlled and controller elements so that their union results in a satisfactory total system. In many cases the controlled element is more or less predetermined by factors beyond the scope of the system designer. In other situations, it is assumed to be known in order to facilitate the design process. The controlled element then is to be regarded as unalterable, and will often be referred to as the "unalterable element." Sensing and actuating elements are "quasi-alterable" in practice, since they are capable of change only by selection of a different item of the same general class. Equalization elements are completely alterable within the realm of physical realizability and practicality. With this basis, system synthesis may be defined as the process of determining the properties of a mechanism required to control an unalterable physical element in some desired fashion.

(b) SYSTEM DESIGN PROCEDURE

This subsection outlines general procedures and methods used to synthesize a controller. It will indicate the various interrelated tasks which must be accomplished before a satisfactory controller is developed. This outline of procedure is not to be construed as the only possible way to achieve the desired result. Rather it

is a method that has been found to be most efficient over an extended period of time, and is in a continual process of development. The basic premise is that proper and efficient design must be firmly founded upon physical understanding of system and component characteristics. A direct corollary is that physical understanding is greatly enhanced by extensive use of mathematical models of the system components.

The aim of servo system design is to integrate components into a functional system. In achieving this end, the designer must

1. Establish system requirements.
2. Synthesize a system meeting the requirements, i.e., select and integrate components into the functional system. In this process, the designer must
 - a. "Live with" the unalterable elements.
 - b. Select or design the best quasi-alterable elements (sensors and actuators) available.
 - c. Design equalization to tie the unalterable and quasi-alterable elements into a well integrated functional system meeting the system requirements.

With this background, it is now possible to discuss the various steps of controller design, remembering that they are chronological in a general sense only; there must be considerable feedback and interrelation existing between steps.

1. SPECIFICATION OF:
 - a. Purpose of system.
 - b. General system requirements.

Discussion: The requirements and purposes of a system are partially derivable from operational conditions imposed upon the equipment. Certain other requirements are somewhat adjustable, and are usually set forth as design objectives. Still others are evolved during the design process.

2. DETERMINATION OF UNALTERABLE ELEMENT (CONTROLLED ELEMENT) CHARACTERISTICS.

Discussion: The defining characteristics and operational features of the unalterable element should be thoroughly understood and completely defined. The physical quantities to be controlled must be identified, and the means through which control can be imposed must be established.

3. DETERMINATION OF BASIC FUNCTIONAL BLOCK DIAGRAM.

Discussion: An intimate knowledge of the system requirements and the unalterable element characteristics,

coupled with a detailed understanding of possible means of achieving the ends required, are the main bases for selecting and evolving the proper functional diagram. It is at this point that experience and understanding of physical systems pays off most heavily. In this phase a type of control is established. The ability to measure the controlled variables, and to control the unalterable elements accordingly must be carefully considered. Preliminary equalization is established by the manner in which connections between elements are made.

4. SELECTION OF ACTUATING AND SENSING ELEMENTS

Discussion: In the functional block diagram phase the type of control has been determined, at least as to the general types of sensors to be used. The unalterable elements and system requirements largely determine the characteristics required of the actuating elements; this narrows the field of possible actuators down to a small group of units available, (or capable of development in an allowable time span). Therefore, a very few versions of sensors and actuators need be considered further. It is the selection of these choice few that is desired in this phase. Final selection is made after a considerable portion of the next phase is completed for all likely combinations.

5. DETAILED SYSTEM STUDY

- a. Studies using "normal" unalterable element characteristics.
- b. Studies using critical abnormal unalterable element characteristics.

Discussion: These studies are performed by the use of one or all of the analysis techniques of chapters III and VIII as tools for use in experimenting with the mathematical model. The modifying characteristics of the remaining equalization are found by judicious use of trial and error experimentation with these models. Hence, the study phase usually consists of many detailed computations using the root-locus, open loop—closed loop logarithmic, and analog computer methods or other means which may be available. In the latter portions of the study, as much information about the physical system is used as is possible, including all important non-linearities. All assumptions are carefully listed in great detail for later verification by actual tests.

6. SYSTEM DETAIL DESIGN

Discussion: This phase consists of taking the mathematical models of the equalization, sensing, and actuating element and designing, developing, and constructing the physical manifestations of these models.

7. TESTS

Discussion: This phase usually starts with individual component tests (which sometimes commence in phase 4 or 5), and ends with tests on the complete system. The major aim is to ascertain that the actual physical equipment will perform in accordance with the system requirements and purposes.

8. PRODUCTION DESIGN, QUALIFICATION AND FUNCTIONAL TESTS, ETC.

Discussion: This phase is concerned with all of the items necessary to make a workable prototype system into a suitable production design. Organizationally, the people charged with this responsibility are usually quite removed from the original system designers. However, the entire effort is directed toward production, and this phase is as important in the overall system synthesis as any other.

The above procedure for designing a controller also covers the complete field of servo system synthesis. However, this chapter is concerned only with aspects basic to items 1, 3, and 5.

Restricting the discussion to basic items then limits the subject of synthesis, as understood here, to that of designing suitable equalization for a system composed of more or less unalterable sensing, actuating and controlled elements.

(c) EQUALIZER SYNTHESIS

The analysis problem of chapter III was particularized from a very general one to the problem of finding the closed-loop transfer function from a knowledge of the open-loop transfer function. In this chapter, certain desired properties of the closed-loop system are known, as well as the characteristics of all of the elements of the open-loop, with the exception of the equalization. The problem then becomes one of determining methods of utilizing elements with essentially fixed characteristics in order to achieve a type of behavior that is in no way inherent in the element itself, by interconnecting it with other elements. The third section of this chapter will point out some ways in which this can be accomplished.

However, an objective must be clearly stated before this can be attempted. This objective is expressed in terms of the desired characteristics of the completed system. Therefore, the second section of the chapter deals with the methods of specifying servo system performance.

SECTION 2 — SPECIFICATION OF SERVO SYSTEM PERFORMANCE

The most complete specification of a linear system would be in the form of the desired system transfer functions. However, such specification is usually impractical, and would not even be useful in some cases. Therefore, most systems are specified in terms of quantities putting bounds upon some of the parameters occurring in the possible transfer functions, i.e., quantities which partially define the desired transfer function. Such specifications allow the system

designer leeway in handling the less important features of the system transfer function and emphasize the important features. This section will discuss in detail the various quantities normally used in specifying linear system performance. All of these quantities will be correlated with the properties of the transfer function which they describe.

Quantities of direct interest in specifying servo per-

formance are:

1. Degree of stability.
2. Response time to representative inputs.
3. Accuracy of control.

A servo system is almost always required to be stable. This immediately specifies that the closed-loop transfer function poles be in the left half of the s -plane only. In addition, a reasonable degree of stability is usually required, making it necessary that the dominant closed-loop second order poles have reasonable damping ratios. Since application of the open loop - closed loop logarithmic method and the root locus methods both lead to closed loop pole and zero values explicitly, the stability and degree of stability is always known throughout any analysis problem. By the use of these methods, any desired damping ratio of a dominant closed loop mode can be easily selected directly if desirable. Therefore, the degree of stability requirement could be set up in terms of the dominant closed loop mode damping ratio if such a dominant mode exists alone. If more than one dominant oscillatory mode exists in the system, a minimum damping ratio requirement is often satisfactory.

The response time requirement basically fixes the damping times of the dominant modes, and hence, either the time constants of dominant closed-loop modes or the damping ratio-undamped natural frequency products, or both. This information is also directly available from the use of the analytical methods previously discussed.

The steady state error of a closed-loop system is a matter of great interest in many applications, and has often been used as a defining characteristic of servo systems. Consider, for example, the servo system of figure IV-1.



Figure IV-1. Illustrative Servomechanism

The error-input transfer function is given by:

$$(IV-1) \quad \frac{E(s)}{R(s)} = \frac{1}{1 + Y(s)} = \frac{1}{1 + \frac{KN(s)}{s^n D(s)}} = \frac{s^n D(s)}{s^n D(s) + KN(s)}$$

The steady-state error is given by use of the final value theorem of the Laplace transformation as:

$$(IV-2) \quad \lim_{t \rightarrow \infty} e(t) = \lim_{s \rightarrow 0} s \left[\frac{R(s)}{1 + Y(s)} \right] = \lim_{s \rightarrow 0} \frac{s^{n+1} D(s)}{s^n D(s) + KN(s)} R(s)$$

The limit as $s \rightarrow 0$ of $D(s)$ and $N(s)$ is, of course, equal to 1, since $N(s)$ and $D(s)$ end in 1.

$$(IV-3) \quad \lim_{t \rightarrow \infty} e(t) = \lim_{s \rightarrow 0} \frac{s^{n+1} R(s)}{s^n + K}$$

If $R(s)$ is a unit step displacement, i.e. $R(s) = 1/s$, and $n > 1$, the steady state error is zero. Hence, a servo system having an open-loop transfer function of the form of (IV-4):

$$(IV-4) \quad \frac{C(s)}{E(s)} = \frac{KN(s)}{s^n D(s)}$$

is called a zero position error system.* If $R(s)$ represents an input velocity step, i.e. $R(s) = 1/s^2$, and $n > 2$, the steady state error will be zero for this type of input also. Therefore, systems with open-loop transfer functions of the form of (IV-5):

$$(IV-5) \quad \frac{C(s)}{E(s)} = \frac{KN(s)}{s^n D(s)}$$

are called zero velocity error systems. Similarly, if $n > 3$, the system is a zero acceleration error system, and so forth, ad infinitum.

If $n = 0$, the system of figure IV-1 has a steady-state error to a unit step displacement input of

$$(IV-6) \quad \lim_{t \rightarrow \infty} e(t) = \lim_{s \rightarrow 0} s E(s) = \lim_{s \rightarrow 0} \frac{s D(s)}{D(s) + KN(s)} \cdot \frac{1}{s} = \frac{1}{1 + K}$$

This system error can be expressed as:

$$(IV-7) \quad \lim_{t \rightarrow \infty} e(t) = \frac{1}{1 + K} = C_0$$

(Since $R(t)$ in this case is a step function $\lim_{t \rightarrow \infty} R(t)$ is simply the magnitude of the step function). C_0 is called the position error coefficient. Similarly, if a unit step velocity is applied to a zero position error system, the steady state position error is:

$$(IV-8) \quad \lim_{t \rightarrow \infty} e(t) = \frac{1}{K} \frac{dR(t)}{dt} = C_v \frac{dR(t)}{dt}$$

where $\lim_{t \rightarrow \infty} dR(t)/dt$ is simply the constant velocity of the input; and C_v is called the velocity error coefficient.

The concept of an error coefficient may be generalized, to provide a very useful way of considering the nature of the response of a system to almost any arbitrary input.

Consider the error-input transfer function for a general servo system, as $M(s)$.

$$(IV-9) \quad \frac{E(s)}{R(s)} = M(s)$$

Assume that $M(s)$ can be expanded in a power series in s which is valid for some values of s . Then

$$(IV-10) \quad E(s) = M(s)R(s) = A_0 R(s) + A_1 s R(s) + A_2 s^2 R(s) + \dots$$

The region of convergence of this power series is near $s = 0$. Since $s \rightarrow 0$ is equivalent in the time domain to $t \rightarrow \infty$, an expression for $E(t)$ as $t \rightarrow \infty$ may be obtained from (IV-10). It can be shown, by utilizing the properties of asymptotic behavior of functions, that (IV-10) may be inverse transformed term by term,

* It should be noted that the term 'zero position error' as given here implies that the output has zero steady-state error only if the 'position' input is a constant in the steady-state. If the input is of a different nature, the so-called zero position error system will actually have a steady-state 'position' error.

Chapter IV Section 2

giving

$$(IV-11) \quad \varepsilon(t) \Big|_{t \rightarrow \infty} = A_0 R(t) + A_1 \frac{dR(t)}{dt} + A_2 \frac{d^2 R(t)}{dt^2} + \dots$$

where $A_0 = C_0$, $A_1 = C_v$, $A_2 = C_a/2$ and C_0 , C_v and C_a are the general position, velocity, and acceleration error coefficients, respectively.

To illustrate the general error coefficient concept, consider the simple servo system of figure IV-2.

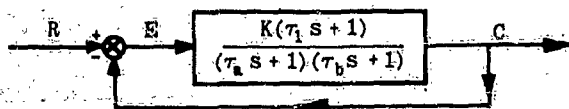


Figure IV-2. Simple Servomechanism

The error-input transfer function is given by:

$$(IV-12) \quad \frac{E(s)}{R(s)} = \frac{1}{1 + \frac{KN(s)}{D(s)}} = \frac{(\tau_a s + 1)(\tau_b s + 1)}{(\tau_a s + 1)(\tau_b s + 1) + K(\tau_1 s + 1)}$$

$$= \frac{[1 + (\tau_a + \tau_b)s + \tau_a \tau_b s^2]}{(1 + K) \left[1 + \frac{\tau_a + \tau_b + K\tau_1}{1 + K} s + \frac{\tau_a \tau_b}{1 + K} s^2 \right]}$$

Since the bracketed term in the denominator is of the form $[1 + Z(s)]$, and since, by the binomial theorem:

$$(IV-13) \quad \frac{1}{1 + Z(s)} = 1 - Z(s) + Z^2(s) - \dots + \dots$$

the expression for the error may be developed into a power series:

$$(IV-14) \quad E(s) = \frac{R(s)}{1 + K} \left\{ \left[1 + (\tau_a + \tau_b)s + \tau_a \tau_b s^2 \right] \left[1 - \left(\frac{\tau_a + \tau_b + K\tau_1}{1 + K} s + \frac{\tau_a \tau_b}{1 + K} s^2 \right) + \left(\frac{\tau_a + \tau_b + K\tau_1}{1 + K} \right)^2 s^2 - \dots \right] \right\}$$

Multiplying and collecting terms; $E(s)$ becomes

$$(IV-15) \quad E(s) = \frac{R(s)}{1 + K} \left\{ \left\{ 1 + \frac{K}{1 + K} [\tau_a + \tau_b + \tau_1] s + \frac{K}{1 + K} \tau_a \tau_b \right. \right. \\ \left. \left. + \frac{1}{1 + K} [(1 - K)(\tau_a + \tau_b)\tau_1 + K\tau_1^2 - (\tau_a + \tau_b)^2] s^2 + \dots \right\} \right\}$$

or

$$(IV-16) \quad E(s) = C_0 R(s) + C_v s R(s) + \frac{C_a}{2} s^2 R(s) + \dots$$

where $C_0 = \frac{1}{1 + K}$, $C_v = \frac{K}{(1 + K)^2} [\tau_a + \tau_b + \tau_1]$, and

$$C_a = \frac{2K}{(1 + K)^2} \left\{ \tau_a \tau_b + \frac{1}{1 + K} [(1 - K)(\tau_a + \tau_b)\tau_1 + K\tau_1^2 - (\tau_a + \tau_b)^2] \right\}$$

Applying the inversion utilized to obtain equation (IV-11), the steady-state error expression becomes:

$$(IV-17) \quad \varepsilon(t) \Big|_{t \rightarrow \infty} = C_0 R(t) + C_v \frac{dR(t)}{dt} + \frac{C_a}{2} \frac{d^2 R(t)}{dt^2} + \dots$$

If $R(t)$ is given by, say,

$$(IV-18) \quad R(t) = vt$$

the steady-state error is approximated by:

$$(IV-19) \quad \varepsilon(t) \Big|_{t \rightarrow \infty} = \frac{1}{1 + K} vt + \frac{K}{(1 + K)^2} [\tau_a + \tau_b + \tau_1] v$$

The error in the above case becomes infinite as time approaches infinity.

It is of interest to note that all of the error coefficients are decreased with an increase in gain. An increase in the open-loop transfer function denominator time constants causes an increase in the velocity error coefficient, while an increase of numerator time constants causes a decrease of the same error coefficient.

In some instances a requirement that the servo be a zero position, velocity, or acceleration error device is imposed. In these cases, the form of the open-loop transfer function near $s = 0$ is, of course, fixed. Accuracy requirements are often more general than this, however, in which cases it is usual practice to specify the first few error coefficients. The specification of error coefficients immediately puts several bounds upon the closed-loop transfer function in addition to those established by the stability and response time requirement.

Another specification occasionally imposed upon a system is concerned with its response to disturbances and "unwanted inputs." If the unwanted input is random, there are methods of determining the root-mean-square error due to the disturbance. These methods are considered in detail in Chapter V. For the present, it is sufficient to note that the specification of an allowable rms error due to unwanted inputs also places bounds upon the closed-loop transfer function.

In view of the above discussion, it is clear that specifications in terms of actual operating conditions may be set up for a given system, and that these requirements define certain properties of the complete system transfer function. If the system unalterable elements are then identified, it becomes possible to consider the methods of connecting some other elements such that the characteristics of the combined system meet the specified system objectives. The next section discusses the basic forms of modification which are used in actual system synthesis.

SECTION 3 — EQUALIZATION

Section IV-1 has defined equalization as the process of modifying performance of elements by external means. Section IV-2 has discussed some of the requirements which would be imposed by specifications

on the performance of a servomechanism. This section will discuss the methods of equalization available for modifying the performance of elements or systems.

It must be recognized that there are an infinite number of ways in which an element or system can be modified. No attempt will be made here to categorize all possible methods of modification. Instead, the discussion will emphasize some of the basic underlying operations that are physically possible, and the fundamental physical forms which can be obtained by use of these operations. Synthesis then consists of properly combining these forms so that the combination meets certain specific requirements. This is usually a cut and try process in which characteristics are assumed for certain available elements, the combination analyzed, changes made, another analysis performed, etc., until satisfactory performance (and hence modification) has been achieved. The trial and error process involved is essentially a series of experiments performed on paper, using mathematical models. An exception to this experimental process is discussed in Chapter V, where "optimum" synthesis procedures are developed which lead directly to a system meeting a specific criterion.

The most basic operational functions that are available for making modifications may be reduced to:

1. Addition (or subtraction)
2. Multiplication (or division)

These operations are denoted graphically in figure IV-3.

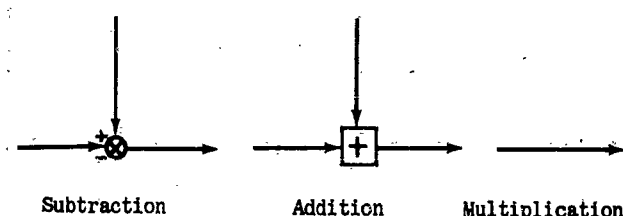


Figure IV-3. Basic Operations

From these operations, three basic structure forms may be derived as shown in figure IV-4.

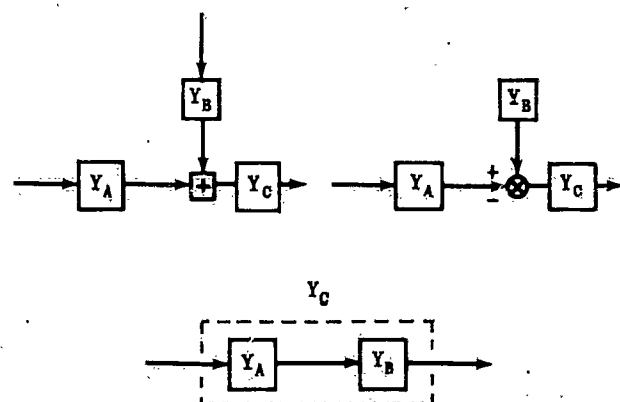


Figure IV-4. Basic Structures

Note that any of the blocks Y_A , Y_B or Y_C may be replaced by a line (or wire), which reduces the block to the trivial case of multiplying by one.

As mentioned above, the basic structure forms shown in figure IV-2 may be used in an indefinitely great number of ways to modify a system or element. There

are, however, several simple combinations of the basic structures which may be considered as fundamental, in that they may be obtained by combining the basic forms using each only once. If multiplication and subtraction are used only once, the single loop feedback combination of figure IV-5 is formed. (Using the addition structure instead of the subtraction leads to essentially the same results.) If multiplication, addition, and subtraction are used only once in a structure, the open loop - closed loop structures of figure IV-6 are obtained.

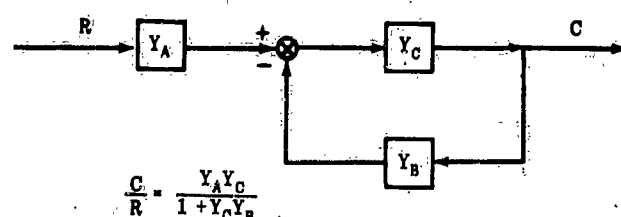
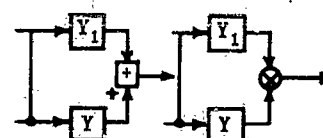


Figure IV-5. Single Loop Feedback Combination

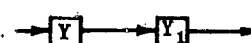
All of the structures, both basic and simple combinations, can be utilized to modify the characteristics of an element or system. These modifications may be logically looked at as transformations. Some simple modifications to an element derivable from the basic structure and first simple combination are given below. Y is the basic element, with Y_1 and Y_2 being modifying (or equalizing) elements.

1. From Basic Structural Forms:

a. $Y \rightarrow Y \pm Y_1$

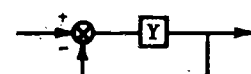


b. $Y \rightarrow Y Y_1$

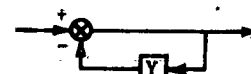


2. From Single Feedback Loop Structure:

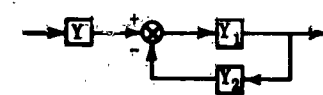
a. $Y \rightarrow \frac{Y}{1+Y}$



b. $Y \rightarrow \frac{1}{1+Y}$



c. $Y \rightarrow Y \frac{Y_1}{1 + Y_1 Y_2}$

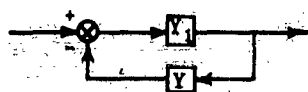


d. $Y \rightarrow Y \frac{1}{1+Y_1}$



Chapter IV Section 3

e. $Y = \frac{Y_1}{1 + Y_1 Y_2}$



h. $Y = \frac{Y_1}{1 + Y_1 Y_2}$



f. $Y = \frac{Y_1 Y_2}{1 + Y_1 Y_2}$



g. $Y = \frac{Y}{1 + Y_1 Y_2}$



The open loop - closed loop structures lead to many more transform pairs.

In synthesis work, an element Y is to be modified to achieve some particular specified form. Experience and knowledge of basic forms usually lead to the proper connections, with trial and error mathematical model experiments leading to the characteristics of the modifying elements.

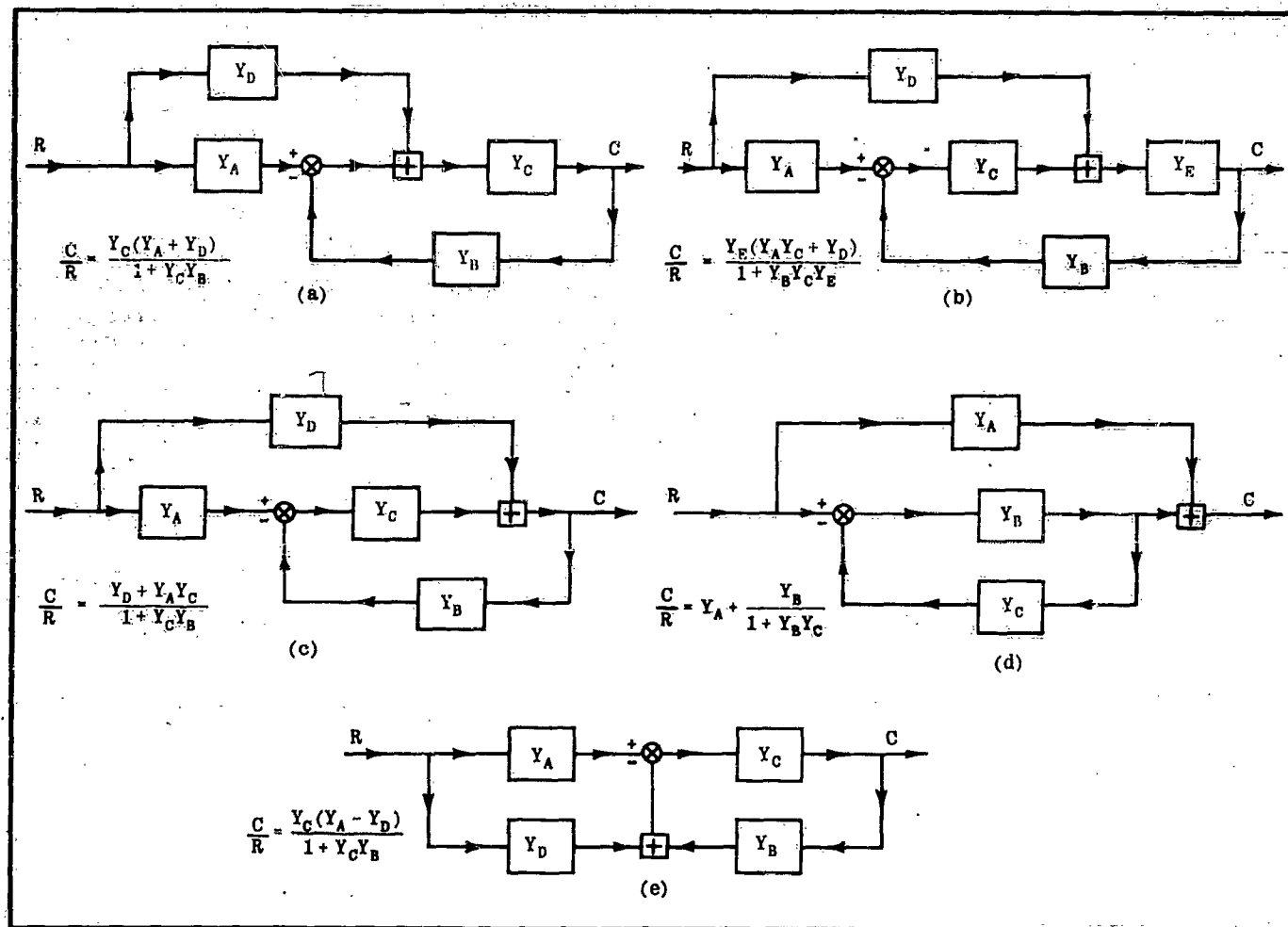


Figure IV-6. Open Loop - Closed Loop Structures

BIBLIOGRAPHY

The following bibliography is included for reference. The list is in no sense complete, but contains the major source material for this chapter. Many of the references, themselves, contain much more complete and detailed bibliographies.

1. 'Theory of Servomechanisms,' by James, Nichols and Phillips; McGraw Hill Book Co., 1947.
2. 'Servomechanisms and Regulating Systems Design,' by H. Chestnut and R. Mayer; John Wiley and Sons, New York, 1951.
3. 'Principles of Servomechanisms,' by G. S. Brown and D. P. Campbell; John Wiley and Sons, New York, 1948.

4. 'Combination Open-Cycle, Closed-Cycle Systems,' by J. R. Moore; Proceedings of the IRE, Vol. 39, 1951.

CHAPTER V

OPTIMUM SYNTHESIS METHODS

Provision is made, at this point, for the addition of a supplementary chapter on "Optimum Synthesis Methods."

While the contents of this volume are essentially complete without the addition of this Chapter V, it is felt that the material it is to contain is of great enough

importance to warrant its inclusion in the interest of making this volume as valuable as possible to the control systems designer.

It is planned to issue this additional chapter in the immediate future.

CHAPTER VI NON-LINEARITIES

SECTION 1 - INTRODUCTION

In most of the preceding work, the assumption has been made that physical systems can be described by linear differential equations of the type

$$(VI-1) \quad \frac{d^n x}{dt^n} + a_1 \frac{d^{n-1} x}{dt^{n-1}} + a_2 \frac{d^{n-2} x}{dt^{n-2}} + \dots + a_{n-1} \frac{dx}{dt} + a_n x = Q(t)$$

where the a_i are constants. This assumption has permitted the development of analytical and graphical methods described in chapters III and IV for the analysis and synthesis of dynamical systems. The validity of this basic assumption was discussed briefly in chapter II.

It is recognized in this chapter that (VI-1) only approximates true physical systems. The coefficients are not ideally constant, but vary as some functions of the dependent variable x . That is, (VI-1) will take the form

$$(VI-2) \quad f_0(x) \frac{d^n x}{dt^n} + f_1(x) \frac{d^{n-1} x}{dt^{n-1}} + \dots + f_n(x) x = Q(t)$$

where the $f_i(x)$ may be any functions of x and its derivatives.

When the coefficients of differential equations are some continuous functions of the dependent variable x , the equations are non-linear differential equations, and the systems described by such equations are said to possess continuous non-linearities. When the coefficients of differential equations are some functions of the variable x , but have finite discontinuities, the equations are non-linear, and the systems are said to have non-linearities of the discontinuous type. The

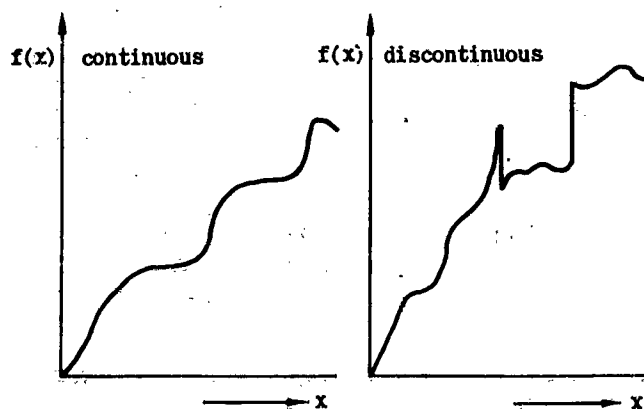


Figure VI-1. Non-Linearities.

two types of non-linearities are shown in Figure (VI-1).

Unfortunately, no general method of analysis or synthesis, such as discussed in chapters III and IV, has been developed for real physical systems having non-linearities. A few writers have developed direct methods of treating some non-linearities, but these methods are generally limited to specific cases and cannot be easily extended to solve the more general problems. However, since an assumed linear system lends itself so readily to analysis and synthesis, an indirect method for the consideration of non-linearities is available. This indirect method consists of finding an answer to the question: "An analysis or synthesis having been made of a system based upon linear constant coefficient equations, what would be the effect of certain non-linearities on the results?" If this question can be answered satisfactorily, the indirect method provides essentially the same information as would a direct solution of the non-linear equation.

The best method for determining the effects of non-linearities on a particular system depends upon the nature of the non-linearities. For this reason, the non-linearities are, as mentioned before, divided into two basic types: continuous and discontinuous. Also since there is a great difference between the effect of large and small discontinuous non-linearities on a system, the discontinuous types are to be further subdivided into major and minor classes.

Continuous non-linearities are usually eliminated from the differential equations of motion by assuming restricted ranges for the variables. By so doing, convenient linear relationships are obtained between such quantities as forces, torques, and moments. It is then desirable to determine if the non-linear terms in the original equations could lead to instability of the system. The first section of this chapter discusses a method, utilizing a theorem due to Liapounoff, by which it is possible to check such continuous non-linear equations for stability.

Discontinuous type non-linearities are unavoidable in real physical systems. They are primarily due to the existence of friction forces, limiting, and free play or hysteresis effects. In addition to these unavoidable discontinuities, some types are added to a physical system to create special effects (spring preloading); or exist by the very nature of the system (relay controllers). These discontinuities are illustrated at the beginning of the second section of this chapter.

The discontinuous type of non-linearity, as such, precludes the application of Liapounoff's theorem, since a basic requirement of the theorem is that the functions $f_i(x)$ in (VI-2) possess continuous derivatives. As mentioned previously, the effects of discontinuous elements depend greatly upon their relative magnitudes. For sufficiently small discontinuities, the physical system approaches the assumed linear system with the discontinuous non-linearities having only a small effect on the system output. This small effect cannot be ignored, however, for its presence may be sufficient to cause sustained oscillations.

If a sinusoidal function is used as the input to an element representing a small discontinuous non-linearity, the output is almost sinusoidal. This similarity to a sinusoid suggests it can be adequately represented by the fundamental of a Fourier expansion. If the discontinuous non-linearity is included in a closed loop and if the transmission of higher order harmonics about the loop is small relative to the transmission of the fundamental, the harmonics may be ignored, and the non-linearity represented by an equivalent transfer function with an amplitude-phase characteristic. In effect, then, the discontinuous system is approximated by a linear system as was done with the continuous non-linearities. With such an aid, a Bode diagram made for an assumed linear system may be

revised, and the effect of the discontinuity determined.

The methods used for determining the equivalent transfer functions of the minor discontinuities are discussed in the second section of this chapter. The possibility of steady-state oscillations due to these non-linearities in systems is also considered.

The third section of this chapter considers the major discontinuous non-linearities. These large discontinuities possess neither continuous derivatives nor negligible harmonics in the sinusoidal responses; hence, the methods discussed to this point cannot be used. Although the analog computer (chapter VIII) may be used for the consideration of such non-linearities (as well as the other types), an analytical or graphical approach is desirable for a better understanding of the system and a check of computer results.

The indirect method discussed in this section involves graphical analysis of simple second order systems which contain one or more major non-linearities. It is recognized that most physical systems are more complex; but, if the general effects of major non-linearities can be found for simple systems, an insight is gained as to possible effects of these major discontinuities on a complex system.

SECTION 2 - CONTINUOUS NON-LINEARITIES

In designing a system based upon the linear approximation of the system differential equations, it is recognized that the motions of the true system following some small disturbance will not be exactly as predicted. The validity of this linear approximation is of particular interest when used to predict the stability of the system. If the solution of the continuous non-linear equations indicates a condition of instability not revealed by the "linearized" form, linear approximations and the methods described in chapters III and IV can not be realistically applied.

The possibility of obtaining incorrect answers to the question of stability from the equations of the linear approximation was investigated by M.A. Liapounoff.* The results of his investigation may be summarized in the following theorem: "If the real parts of the roots of the characteristic equation corresponding to the differential equations of the first approximation are different from zero, the equations of the first approximation always give a correct answer to the question of stability of a non-linear system."*** The theorem assumes only that the non-linear terms of the differential equation may be expanded in a Taylor's series about the equilibrium point in question.

According to this theorem, if all the roots of the linear approximation of the differential equation are negative,

the non-linear system is stable about the point in question, and any small temporary disturbance in the input will result in a temporary disturbance in the output. If, however, any of the roots of the linear approximation of the differential equation are positive, the non-linear system is unstable about the equilibrium point and any small temporary disturbance at the input will result in an output which will diverge from this unstable point.

If any of the roots of the linear approximation of the differential equation about the equilibrium point are zero, the theorem may not be used, and higher order terms of the Taylor series must be considered. Zero roots may result in a "conditionally stable" situation which would depend upon the direction of the disturbance. In servo work, conditionally stable situations are usually as undesirable as absolutely unstable situations and hence, the fact that the theorem does not apply is of little consequence.

It may be pointed out that although the linear approximations of the differential equations indicate stable systems for all amplitudes of disturbances, the theorem applies only for small disturbances.

Figure VI-2 illustrates simple stable and unstable equilibrium positions. If a ball is located at (a), a small disturbance in either direction will result in the return of the ball to (a); and the system represented is stable about the equilibrium point (a). If, however, the ball is located at either (b) or (c) and disturbed, it will leave these equilibrium points, and the system represented is unstable about these points.

* Liapounoff, M.A.; *Problème général de la stabilité du mouvement*; Annals of Mathematical Studies, Vol. 17; Princeton Press.

** Minorsky, N., *Introduction to Non-Linear Mechanics*; J.W. Edwards, 1947; Ann Arbor.

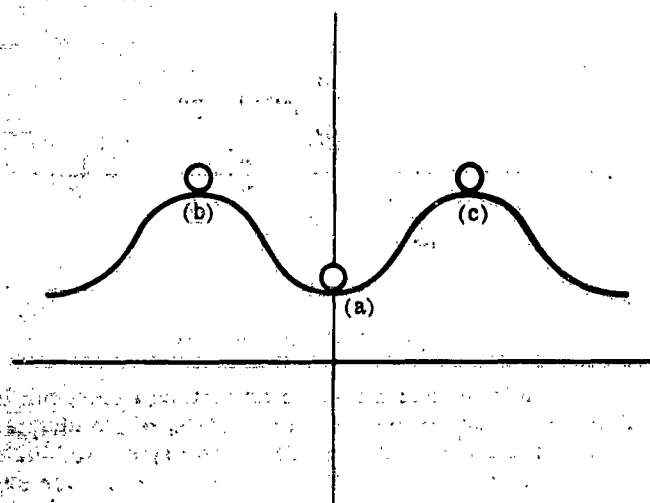


Figure VI-2. Stable and Unstable Equilibrium

To illustrate the application of the theorem, two similar second order non-linear differential equations will be considered:

$$(VI-3) \quad m \frac{d^2x}{dt^2} + \mu(1-x^2) \frac{dx}{dt} + Kx = Q$$

$$(VI-4) \quad m \frac{d^2x}{dt^2} - \mu(1-x^2) \frac{dx}{dt} + Kx = Q$$

When the acceleration and velocity are zero in these equations, the value of x defines the point of equilibrium, x_0 . That is, $Kx = Q$ or $x = Q/K = x_0$.

If a change of variable is made, such as $x = x_0 + \delta$, where δ is the deviation about the equilibrium point x_0 , (VI-3) becomes

$$(VI-5) \quad m \frac{d^2\delta}{dt^2} + \mu[1-(x_0^2 + 2x_0\delta + \delta^2)] \frac{d\delta}{dt} + Kx_0 + K\delta = Q$$

In accordance with the theorem, the first (linear) approximation of the term $\mu[1-(x_0^2 + 2x_0\delta + \delta^2)] \frac{d\delta}{dt}$ is

substituted into equation (VI-5) so that

$$(VI-6) \quad m \frac{d^2\delta}{dt^2} + \mu(1-x_0^2) \frac{d\delta}{dt} + K\delta = 0$$

Equation (VI-6) indicates that the system is stable about all equilibrium points less than unity. As x_0 approaches unity, the coefficient of $d\delta/dt$ becomes very small and the system becomes poorly damped; but still stable for sufficiently small values of δ about x_0 . For $x_0 = 1$, the coefficient of δ becomes zero and one of the roots of the characteristic equation corresponding to (VI-6) will be zero. The theorem is not applicable in this case. For x_0 greater than unity the system described by (VI-6) is unstable.

If equation (VI-4) is analyzed in the same manner, the equation corresponding to (VI-6) is

$$(VI-7) \quad m \ddot{\delta} - \mu(1-x_0^2) \dot{\delta} + K\delta = 0$$

For this case, the coefficient of δ is negative for $0 \leq x < 1$ and the system is unstable about any equilibrium position in this range. For $x_0 > 1$, the coefficient of the δ term is positive, and the system described by (VI-4) is stable for small δ about x_0 .

In general, it may be concluded that the analysis and synthesis of dynamic systems based upon the linear approximations of the continuous non-linear differential equations may be made with assurance that the question of stability will be answered correctly. It must be kept in mind, however, that the answers may be valid only for small disturbances about some equilibrium point. If the response of some dynamical system to a disturbance should be poorly damped, the system should be carefully analyzed; for the presence of some non-linear term in the damping coefficient may cause the true non-linear system to be unstable for very small magnitudes of disturbance. Or, stated in another way, the designer of a system described by linearized equations should satisfy himself that the damping coefficient is relatively independent of the non-linearity.

SECTION 3 - DISCONTINUOUS NON-LINEARITIES

(a) GENERAL

Section VI-1 has dealt with the continuous type of non-linearity. Under certain restrictions, it was found that continuous non-linearities could be linearized, so that the methods of analysis and synthesis outlined in previous chapters could be used. Unfortunately, another type of non-linearity also occurs in most real physical systems. This type is the discontinuous non-linearity.

The form in which discontinuities occur is varied. Figures VI-3 through VI-7 illustrate the static transfer characteristics of several of the more common types.

Coulomb friction exists in all physical systems in which there is relative motion between two contacting surfaces. The friction force is of constant magnitude and is always in such a direction as to resist the relative motion. The discontinuity in systems with coulomb

friction occurs at the instant of a reversal in relative motion. Coulomb friction may be ignored only if the friction force is much smaller than the other forces acting.

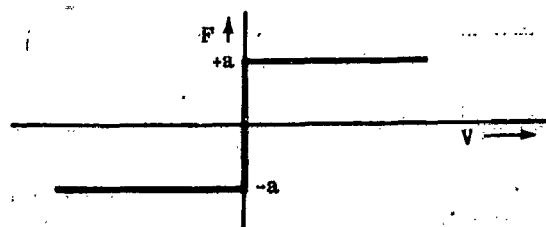


Figure VI-3. Coulomb Friction

In certain cases, it is desirable to spring load a device so that, when external forces are removed, the spring will cause the device to seek a null or zero position. The presence of coulomb friction, however, may cause

the device to stop short of the desired zero position. To overcome this, the spring may be preloaded by an amount equal to the coulomb friction force. Preloading introduces a discontinuity, but is desirable for the above reason.

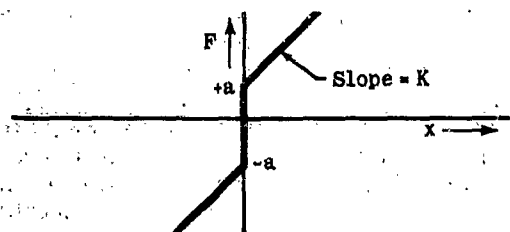


Figure VI-4. Spring Preload

In many physical systems, an input to the system must exceed a certain minimum value before any output is realized. One effect of this type is known as threshold or flatspot and the value which must be exceeded is referred to as the threshold value. Threshold is generally undesirable, and can be ignored only if the threshold value is much smaller than the input value.

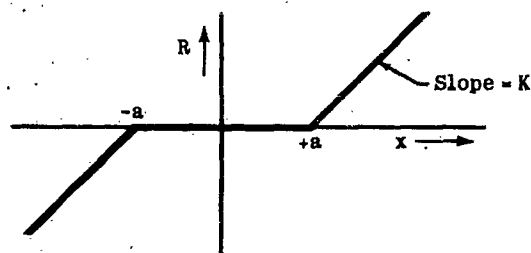


Figure VI-5. Threshold

Linear systems can transmit signals of infinite magnitude. However, all components forming a real system have limitations on such quantities as position, speed, or voltage. In some cases, the limitations may never be reached, with the result that they need not be considered. In other cases, the limitations may be exceeded and cannot be ignored. Limiting may or may not be desirable.

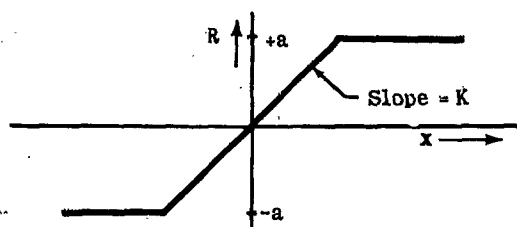


Figure VI-6. Limiting

In many physical devices, a plot of the input versus the output results in a closed curve called a hysteresis loop. The cause of such a loop in mechanical systems is referred to as the backlash or free play which exists between two mechanically coupled components. Hysteresis is generally undesirable and should be eliminated wherever possible.

For the cases in which the discontinuities may not be neglected without serious effect upon the results, it

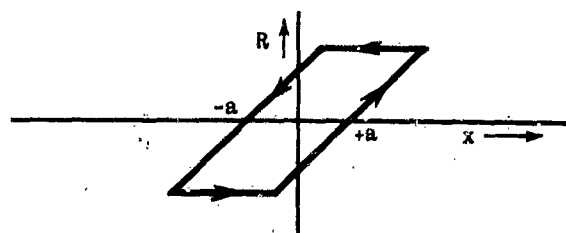


Figure VI-7. Hysteresis

is desirable to have a method or methods available to determine such effects. This section of the chapter will discuss two such methods. The first method is applicable only to small discontinuities; while the second may be used for any discontinuity, providing it occurs in a second or lower order system.

(b) SMALL DISCONTINUITIES

For ease in the analysis and synthesis of closed loop systems, one of the methods of previous chapters made use of Bode diagrams in which the amplitude ratios and phase angles of transfer functions were plotted versus frequency. While these diagrams have not been emphasized as frequency responses, they can be considered as such, since the curves are exactly those which would be obtained if the systems were stable and were excited with sinusoids of varying frequency.

In this subsection, the non-linear elements will be replaced by "equivalent" linear elements. An "equivalent" transfer function will be derived by applying a sinusoidal input to the non-linear element and by determining the Fourier series of the output waveform. The "equivalent" amplitude ratio will be defined as the ratio of the fundamental output amplitude to the input amplitude. The phase angle will be defined as the difference between the phase of the fundamental and that of the input. The concept of such an "equivalent" transfer function is sound if:

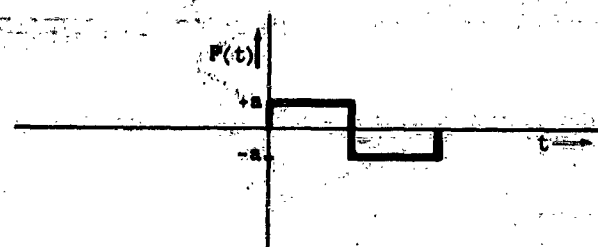
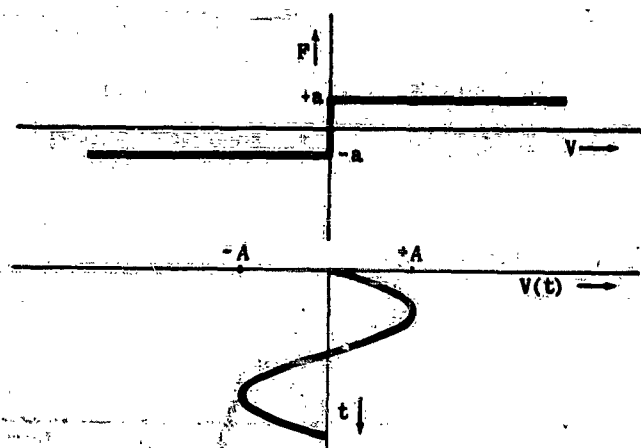
1. The system is oscillating at constant amplitude at the frequency considered. (This implies that the system is either unstable, or that a sinusoidal input is being applied.)
2. The amplitude of harmonics appearing at the input of the non-linear element is negligible. This means that the transmission of these harmonics through the system (around the loop) is negligible compared to the transmission of the fundamental.

A "small discontinuity" may now be defined as one satisfying the above conditions. Since the non-linearities considered here are not frequency sensitive, the ratio between the amplitude of the fundamental of the output wave and that of the input represents a shift in the amplitude ratio of the Bode plot. A similar shift may occur in the phase curve.

The remaining portions of this subsection discuss the equivalent transfer functions of various non-linearities, and their effect upon Bode plots. Ratios of harmonic to input amplitude are derived to aid in determining the validity of condition 2 above for a specific problem.

To obtain the frequency responses of discontinuous elements of the type discussed in this section, it is only necessary to analyze them at a single frequency, for they are not frequency sensitive. In general, the outputs may be represented by Fourier expansions of the responses written as functions of the amplitudes of the discontinuities and the amplitudes of the sine-

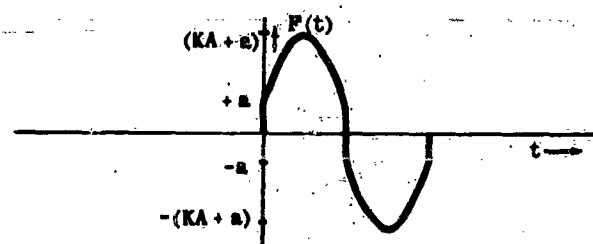
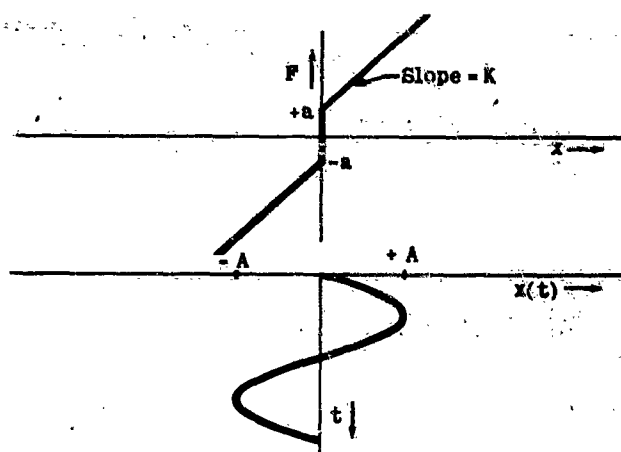
waves being transmitted through the elements. Figures VI-8 through VI-12 illustrate the nature of the responses of discontinuous elements to sinusoidal inputs. If the amplitudes of the harmonics may be ignored when compared to the amplitudes of the fundamentals, the amplitudes and phases of the fundamentals plotted against frequency would be the frequency responses or transfer



Input: $V(t) = A \sin \omega t$

Output: $F(t) = \frac{4a}{\pi} \left[\sin \omega t + \frac{1}{3} \sin 3\omega t + \frac{1}{5} \sin 5\omega t + \dots \right]$

Figure VI-8. Transfer Characteristic of Coulomb Friction



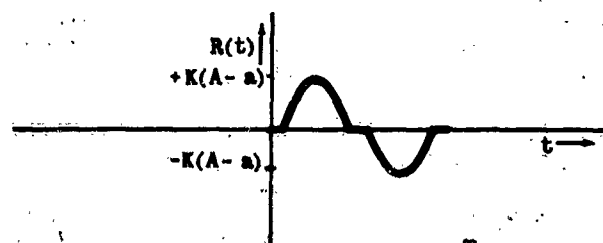
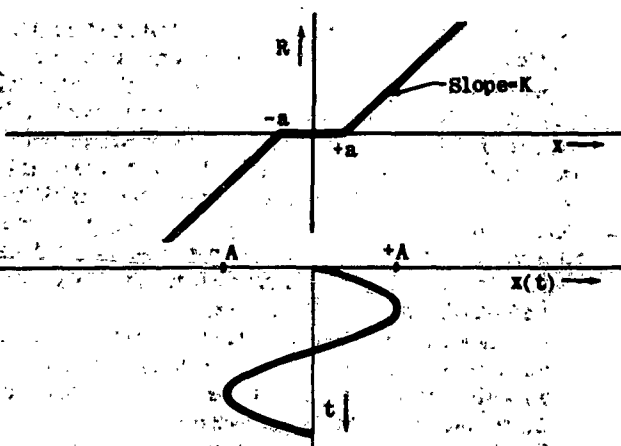
Input: $x(t) = A \sin \omega t$

Where: $b_1 = AK + \frac{4a}{\pi}$

Output: $F(t) = b_1 \sin \omega t + \sum_{n=3,5,\dots} b_n \sin n\omega t$

$b_n = \frac{4a}{n\pi}$

Figure VI-9. Transfer Characteristic of Spring Preload



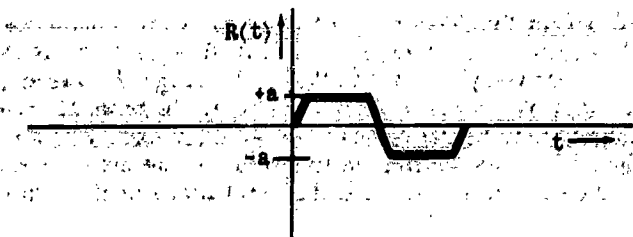
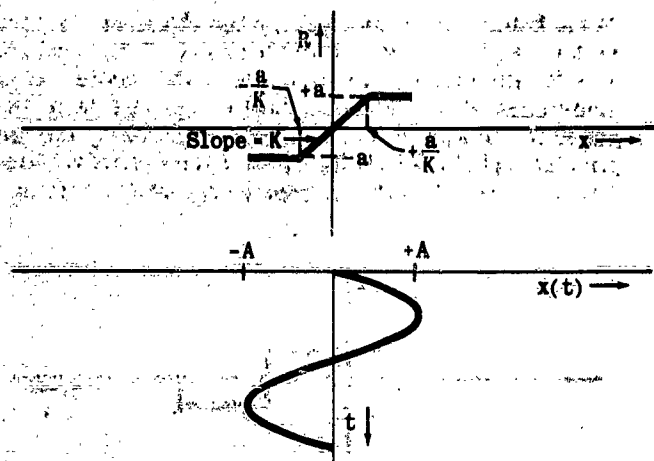
Input: $x(t) = A \sin \omega t$ Output: $R(t) = b_1 \sin \omega t + \sum_{n=3,5,\dots} b_n \sin n\omega t$

Where: $b_1 = \frac{AK}{\pi} [\pi - 2B + \sin 2B] - \frac{4Ka}{\pi} \cos B$

$b_n = \frac{4K}{\pi(1-n^2)} (\sin nB \cos B - n \sin B \cos nB) - \frac{4a}{\pi} \cos nB$

$B = \sin^{-1} \frac{a}{A}$

Figure VI-10. Transfer Characteristic of Threshold



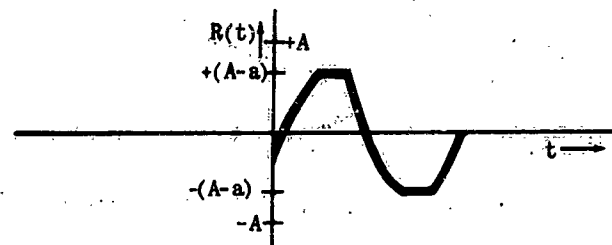
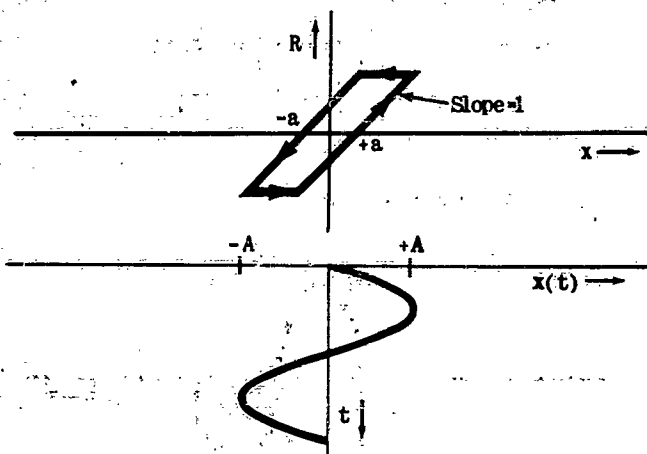
Input: $x(t) = A \sin \omega t$ Output: $R(t) = b_1 \sin \omega t + \sum_{n=3,5,\dots}^{\infty} b_n \sin n \omega t$

$$\text{Where } b_1 = \frac{1}{\pi} \left[AK \left(2 \sin^{-1} \frac{a}{AK} - \frac{2a}{AK} \sqrt{1 - \left(\frac{a}{AK} \right)^2} \right) + 4a \sqrt{1 - \left(\frac{a}{AK} \right)^2} \right]$$

$$b_n = \frac{4}{\pi} \left[\left(\frac{AK}{1-n^2} \right) (n \sin B \cos nB - \sin nB \cos B) + \frac{a}{n} \cos nB \right]$$

$$B = \sin^{-1} \frac{a}{AK}$$

Figure VI-11. Transfer Characteristic of Limiting



Input: $x(t) = A \sin \omega t$ Output: $R(t) = z_1 \sin(\omega t + \phi_1) + \sum_{n=3,5,\dots}^{\infty} z_n \sin(\omega t + \phi_n)$

$$\text{Where: } z_1 = \frac{A}{\pi} \sqrt{1 - u^2 + \left(\frac{3\pi}{2} - \sin^{-1} u \right)^2} + 2 \left(\frac{3\pi}{2} - \sin^{-1} u \right) u \sqrt{1 - u^2}$$

$$\phi_1 = \tan^{-1} \frac{u^2 - 1}{\left[\frac{3\pi}{2} - \sin^{-1} u + u \sqrt{1 - u^2} \right]}$$

$$z_n = \sqrt{a_n^2 + b_n^2}$$

$$\phi_n = \tan^{-1} \frac{a_n}{b_n}$$

$$a_3 = -\frac{A}{6\pi} \left[\cos 2P + \frac{1}{2}(1 + \cos 4P) \right]$$

$$b_3 = -\frac{A}{6\pi} \left[\frac{1}{2} \sin 4P + \sin 2P \right]$$

$$P = \sin^{-1} u$$

$$a_5 = -\frac{A}{10\pi} \left[\frac{1}{8} - \frac{1}{2} \cos 4P - \frac{1}{3} \cos 6P \right]$$

$$b_5 = -\frac{A}{10\pi} \left[\frac{1}{3} \sin 6P + \frac{1}{2} \sin 4P \right]$$

$$u = 1 - 2 \frac{a}{A}$$

Figure VI-12. Transfer Characteristic of Hysteresis

functions of the discontinuities.

In Figures VI-13 through VI-16 the amplitude ratios and phase shifts of the fundamentals and harmonics are plotted against the ratio of discontinuity amplitudes to element input amplitudes (a/A) to show how the relative amplitude ratios vary.*

From these figures, it is apparent that for very small discontinuities, the harmonic amplitudes are negligible compared to the fundamentals. However, as the discontinuities become larger, these figures show an increase in harmonic amplitudes relative to the fundamentals. In these latter cases the system equations must be examined to determine if condition 2 above will hold. This can be done with the aid of the Bode diagram. If the harmonics are not negligible, the discontinuities are considered as large and will not be dis-

cussed in this part of Section 3. It is interesting to note that the discontinuity amplitude is not necessarily small relative to the input amplitude for the discontinuity to be considered small, e.g., the backlash case.

To apply the amplitude ratios and phases of Figures VI-13 through VI-16 to a particular problem, it is necessary to have an open loop Bode diagram for the system. The primary frequency of interest is in the region of the crossover point (the point at which the amplitude ratio curve passes through zero db), because it is within this region that the stability of a closed-loop system is determined. The procedure to determine effects of discontinuities is as follows:

1. Draw the Bode diagram for the linear system and establish an optimum gain.
2. Assume that the system is oscillating at a frequency corresponding to the crossover point (point at which amplitude ratio curve intersects zero db line) and at constant amplitude.
3. Assume that the harmonics of the Fourier series representing the output wave of the non-linear element are negligible.

* The 'amplitude ratio' of a harmonic to a fundamental is an arbitrary term, and completely ignores the difference in frequency existing between the harmonic and the fundamental

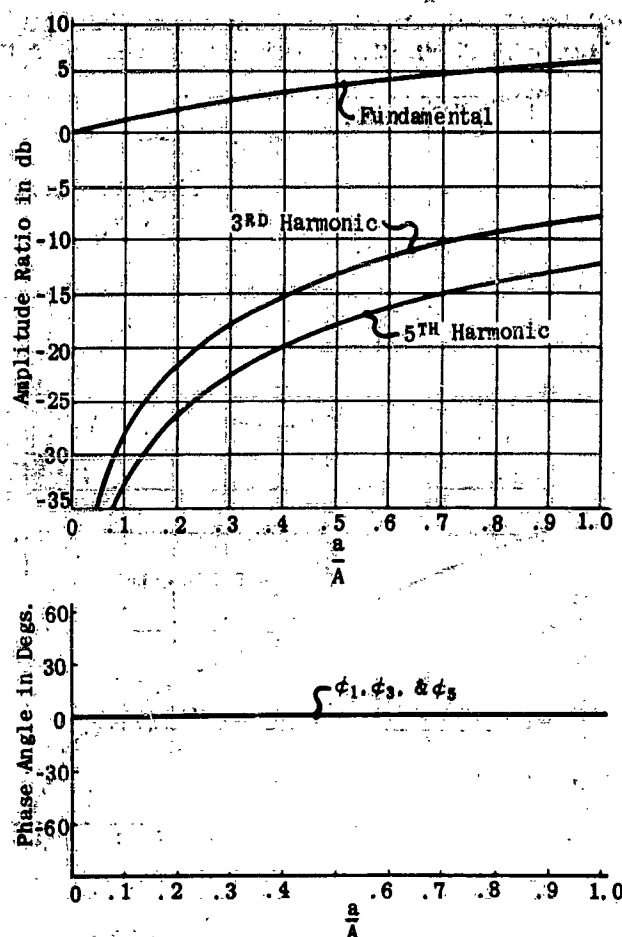


Figure VI-13. Amplitude Ratio and Phase Due to Spring Preload

4. Determine effect that the fundamental of the non-linear element output has on the Bode diagram.
5. Check the validity of step 3.

In general, if the amplitude ratio of the discontinuous element is greater than unity (zero db), the closed-loop system becomes less stable. If the gain is less than zero db, the closed-loop system becomes more stable. As will be shown in one of the examples to follow, it is possible to have a combination of amplitude ratio and phase lag with a resulting instability of a system which is stable when discontinuities are not present. Several examples will now be discussed in which the method outlined above will be applied to the simple positional servo illustrated in figure VI-17.

The closed and open-loop differential equations describing this system are, respectively:

$$(VI-8a) \quad \tau \frac{d^2C}{dt^2} + \frac{dC}{dt} + K_1 K_2 C = K_1 R$$

$$(VI-8b) \quad \frac{d^2B}{dt^2} + \frac{1}{\tau} \frac{dB}{dt} = \frac{K_1 K_2}{\tau} E$$

In each example, it will be initially assumed that the harmonics have a negligible effect. Following this, the validity of the assumption will be checked.

COULOMB FRICTION. Figure VI-8 illustrates the transfer characteristics relating the coulomb friction force to a sinusoidal velocity. When the harmonics are neglected, the response of the non-linearity is given by the first term of the Fourier series, or $F(t) = (4a/\pi) \sin \omega t$. The transfer characteristics corresponding to this first approximation are illustrated in figure VI-18. The figure shows that the slope or "gain" has the units of viscous friction, that is, the first approximation to coulomb friction is viscous friction.

If the output of the closed-loop system of figure VI-17 is subjected to this "effective" viscous friction, the differential equation written with the first approximation becomes

$$(VI-9) \quad \tau \frac{d^2C}{dt^2} + \frac{dC}{dt} \left(1 + \frac{4a}{\pi A}\right) + K_1 K_2 C = K_1 R$$

From this expression, it is concluded that the effect of coulomb friction on the system of figure VI-17 (viewed in this way) is to increase the damping ratio without changing the frequency. It is to be further noted that for a given value of coulomb friction, the effective coefficient of viscosity approaches infinity as the input amplitude approaches zero.

Equation (VI-9) may be written in the form of (VI-10)

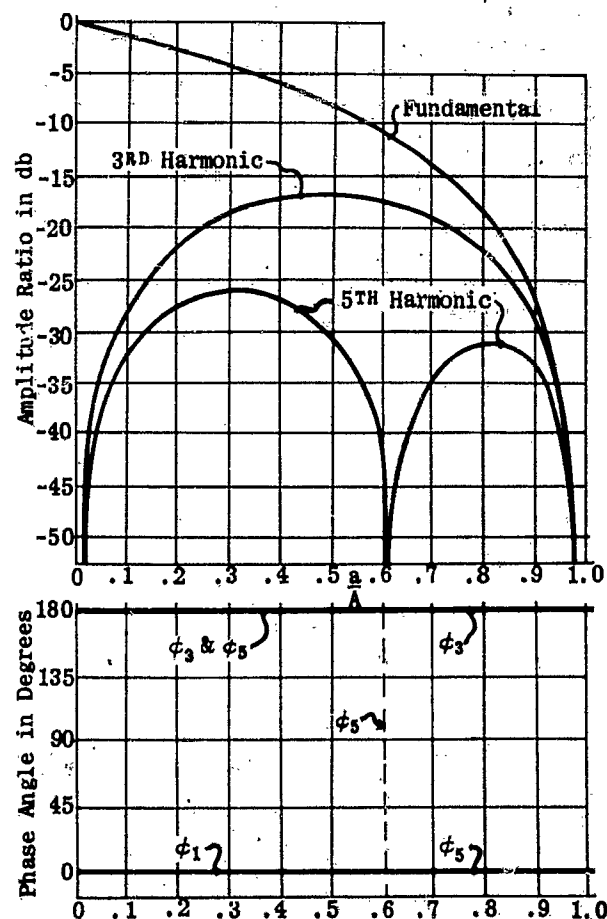


Figure VI-14. Amplitude Ratio and Phase Due to Threshold

Chapter VI
Section 3

$$(VI-10) \quad \tau \frac{d^2C}{dt^2} + \frac{dC}{dt} + K_1 K_2 C = K_1 \left(R_1 - \frac{4a}{\pi K_1 A} \frac{dC}{dt} \right)$$

which, after transforming, becomes

$$(VI-11) \quad C(\tau s^2 + s + K_1 K_2) = K_1 \left(R_1 - \frac{4a}{\pi K_1 A} s C \right)$$

A block diagram corresponding to this latter equation is shown in figure VI-19 which illustrates that the

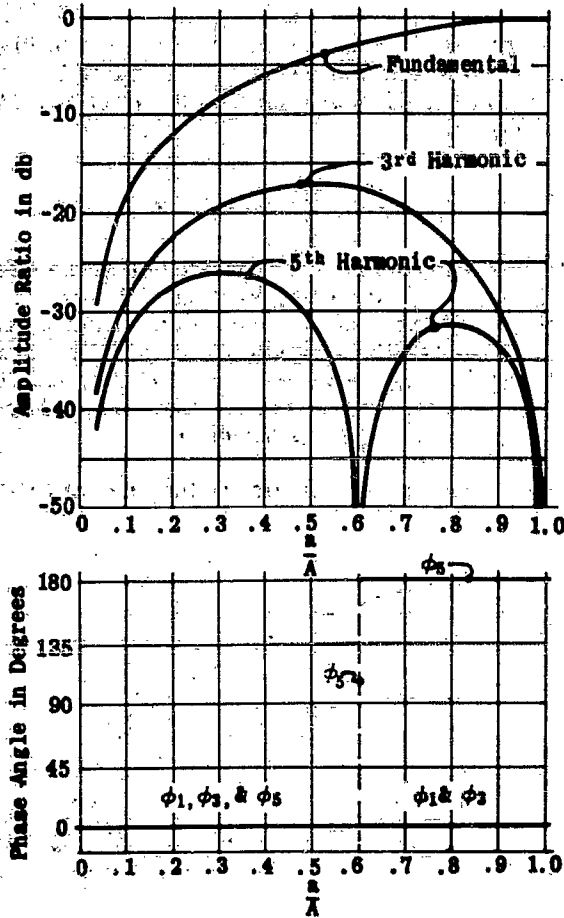


Figure VI-15. Amplitude Ratio and Phase Due to Limiting

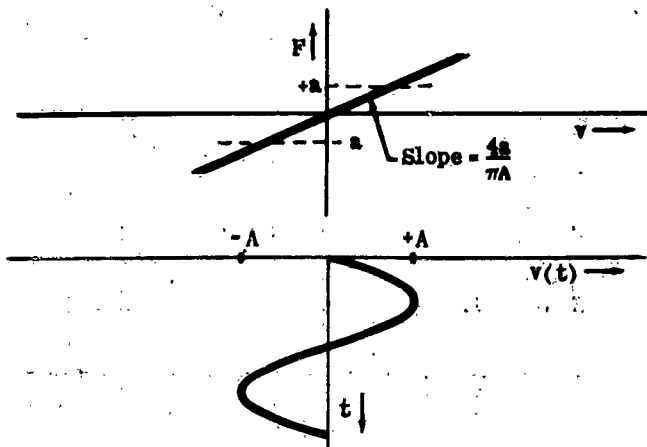


Figure VI-18. Effective Characteristics of Coulomb Friction when Harmonics are Negligible

addition of coulomb friction to the closed-loop system of figure VI-17 effectively adds another feedback path.

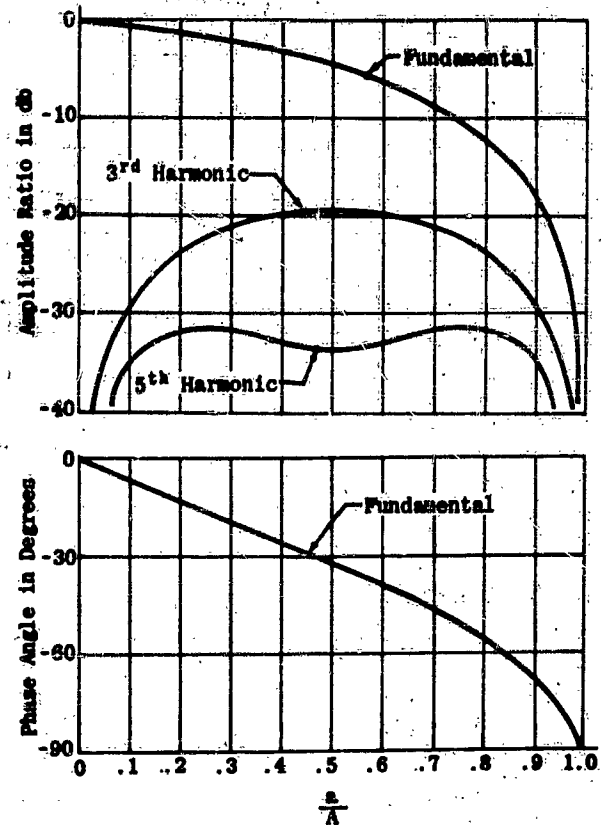


Figure VI-16. Amplitude Ratio and Phase Due to Hysteresis

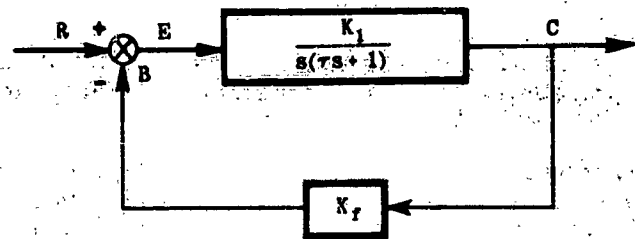
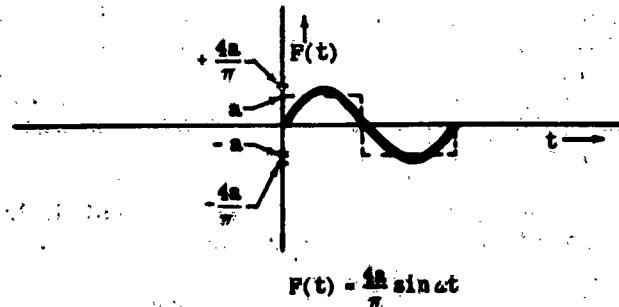


Figure VI-17. Position Servo



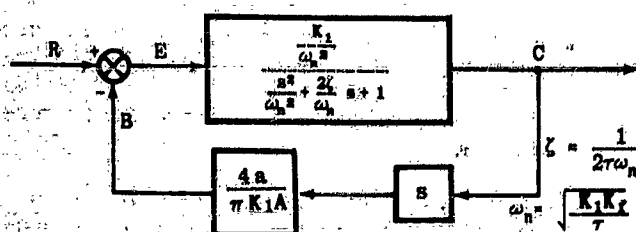


Figure VI-19. Position Servo with Coulomb Friction

The open-loop Bode diagram of figure VI-19 is shown in figure VI-20.

It is evident that decreasing amplitude A of the controlled variable (input to coulomb friction block) increases the open-loop gain, thus resulting in an increase in the apparent damping ratio of the closed-loop system.

The Bode diagram shows that at the frequency of the third harmonic the signal is attenuated 9.5 db below its amplitude ratio at the fundamental frequency. The Fourier expansion accompanying figure VI-8 reveals a further attenuation of three (9.5 decibels) due to the feedback element. With such a large attenuation of the harmonics relative to the fundamental, it can be concluded that the harmonics contribute little to the stability of the system.

SPRING PRELOAD. Figure VI-9 illustrates the transfer characteristics relating the force, applied by a preloaded spring, to a sinusoidal displacement. When the harmonics are neglected, the response of the non-linearity is given by the first term of the Fourier series, or $F(t) = [AK + (4a/\pi)] \sin \omega t$. The transfer characteristics corresponding to this first approximation are illustrated in figure VI-21.

As indicated by this figure, the preloaded spring is effectively replaced with a spring with a coefficient of $K + (4a/\pi A)$. The closed-loop differential equation written for the system of figure VI-17 when loaded with the above spring is

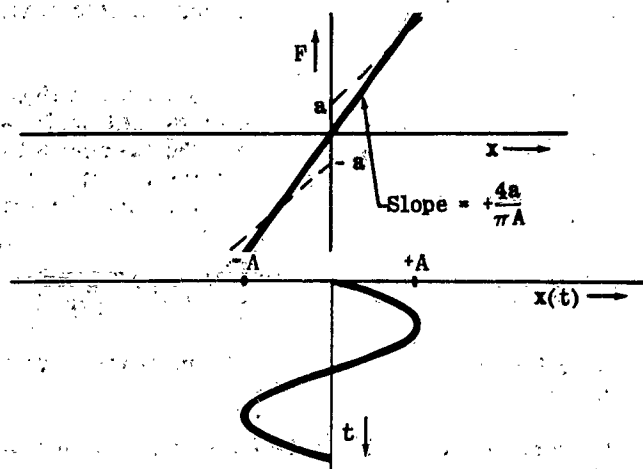


Figure VI-21. Effective Characteristics of Spring Preload when Harmonics are Negligible

$$(VI-12) \quad \tau \frac{d^2 C}{dt^2} + \frac{dC}{dt} + \left[K_1 K_2 + \left(K + \frac{4a}{\pi A} \right) \right] C = K_1 R$$

The addition of this spring increases the natural frequency and decreases the damping ratio of the system.

Writing equation (VI-12) in the form of (VI-13)

$$(VI-13) \quad \tau \frac{d^2 C}{dt^2} + \frac{dC}{dt} + K_1 K_2 C = K_1 \left[R - \frac{\left(K + \frac{4a}{\pi A} \right) C}{K_1} \right]$$

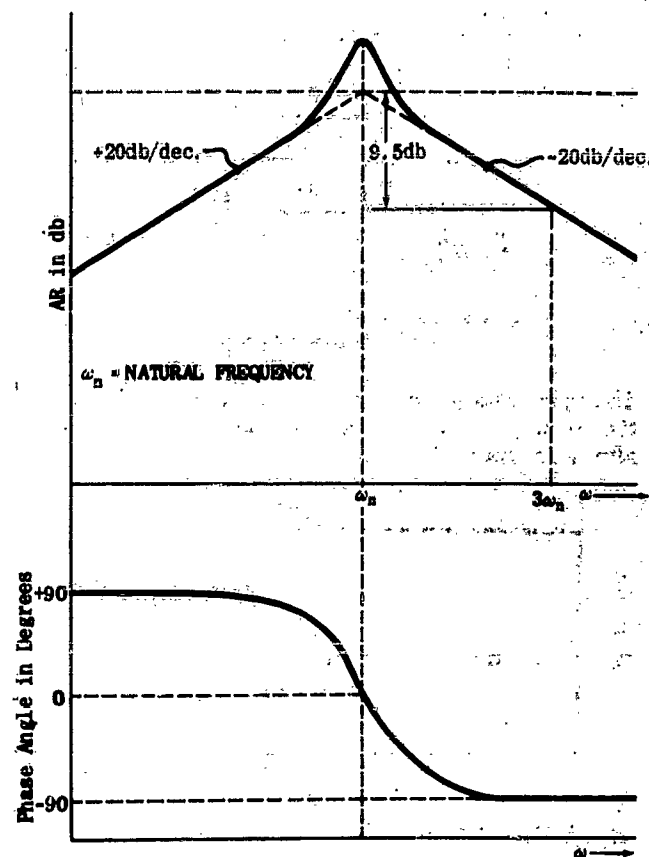
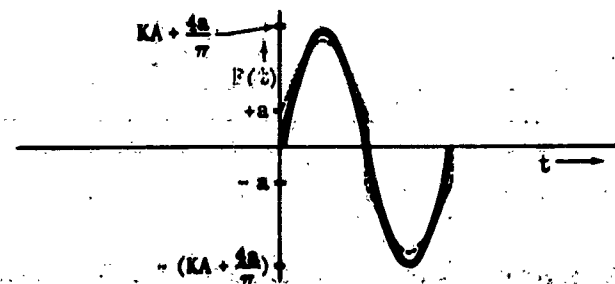


Figure VI-20. Open-Loop Bode Diagram of System with Coulomb Friction



Chapter VI Section 3

permits drawing figure VI-22. As with the previous case, a preloaded spring may be represented by the addition of a second feedback path to figure VI-17.

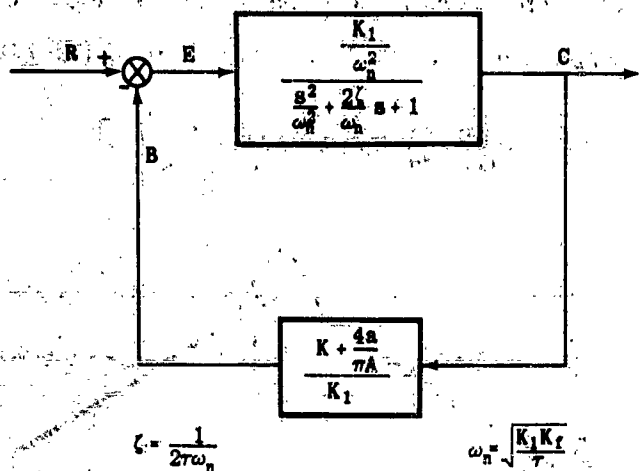


Figure VI-22. Position Servo with Preloaded Spring

The open loop Bode diagram of figure VI-22, in which the term $[K + (4a/\pi A)]/K_1$ is now part of the gain, is shown in figure VI-23.

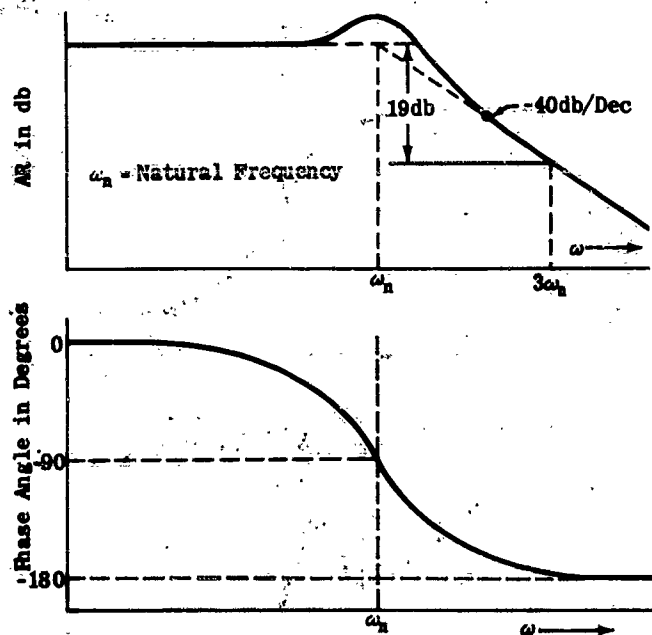


Figure VI-23. Open Loop Bode Diagram of System with Spring Preload

For increasing values of the ratio (a/A) the gain increases with a resulting increase in the natural frequency of the closed loop and a decrease in damping ratio.

As is apparent from the curves of figure VI-13 and the fairly rapid attenuation above the natural frequency shown by the Bode diagram, the harmonics have little effect on the stability as determined by the fundamental alone.

In the above cases it may be noticed that the addition

of non-linear loads to the closed loop system of figure VI-17 effectively added a second feedback path. In the following cases, the feedback path of figure VI-17 will contain the non-linearities and the form of the figure will be unchanged.

THRESHOLD. It will now be assumed that the position servo illustrated in figure VI-17 has threshold non-linearity in the feedback portion of the loop. When the harmonics can be ignored, the curve of figure VI-10 are effectively changed to that of figure VI-24 with a slope of $(K/\pi) (\pi - 2B + \sin 2B) \sim [(4Ka)/(\pi A)] \cos B$ as determined from the first term of the Fourier series.

Figure VI-24 also shows that K_2 becomes part of the gain of the open loop. If K_2 is the effective slope, the differential equation for the system becomes

$$(VI-14) \quad \tau \frac{d^2 C}{dt^2} + \frac{dC}{dt} + K_1 K_2 C = K_1 R$$

From this, it is seen that a decrease in K_2 results in a lower natural frequency and a higher damping ratio of the closed loop system. This is also illustrated in the open loop Bode diagram of figure VI-25.

When the zero db line intersects the -40 db/dec portion of the Bode diagram, the rapid attenuation of the system above the natural frequency suggests the possibility of neglecting the harmonics. For smaller values of inputs to the non-linear element, however, the zero decibel line may cut the -20 db/dec line with the result that the over-all system gain at the third harmonic must be compared with the over-all gain at the fundamental. Although figure VI-25 shows the same harmonic attenuation for both locations of the zero db line, it must be remembered that the straight lines shown are asymptotes to the actual curve. At the natural frequency, ω_n , the asymptotes are below the true curve; and at the natural frequency, ω_n' , the asymptotes may be above the true curve. Therefore, while attenuation in the unprimed case may actually be greater than 19 db, the attenuation in the primed case may be less than 19 db. As shown by the relative gains from figure VI-14, the third harmonic approaches the fundamental in gain as the ratio a/A approaches unity. In this range then, the discontinuity is not small and the stability of the system for small amplitudes of the controlled variable cannot be determined in this way.

LIMITING. The position servo illustrated in figure VI-17 will now be assumed to possess limiting in the feedback portion of the loop. When the harmonics can be ignored, the transfer curve of the non-linearity, figure VI-11, is effectively changed to that of figure VI-26 with a slope

$$K' = \frac{1}{\pi} \left[K \left(\sin^{-1} \frac{a}{AK} - \frac{2a}{AK} \sqrt{1 - \left(\frac{a}{AK} \right)^2} \right) + \frac{4a}{A} \sqrt{1 - \left(\frac{a}{AK} \right)^2} \right]$$

as determined from the first term of the Fourier series.

The differential equation for this example is identical with that for threshold including the discussion which followed. The effect of limiting in the feedback path of figure VI-17 is to reduce the apparent natural fre-

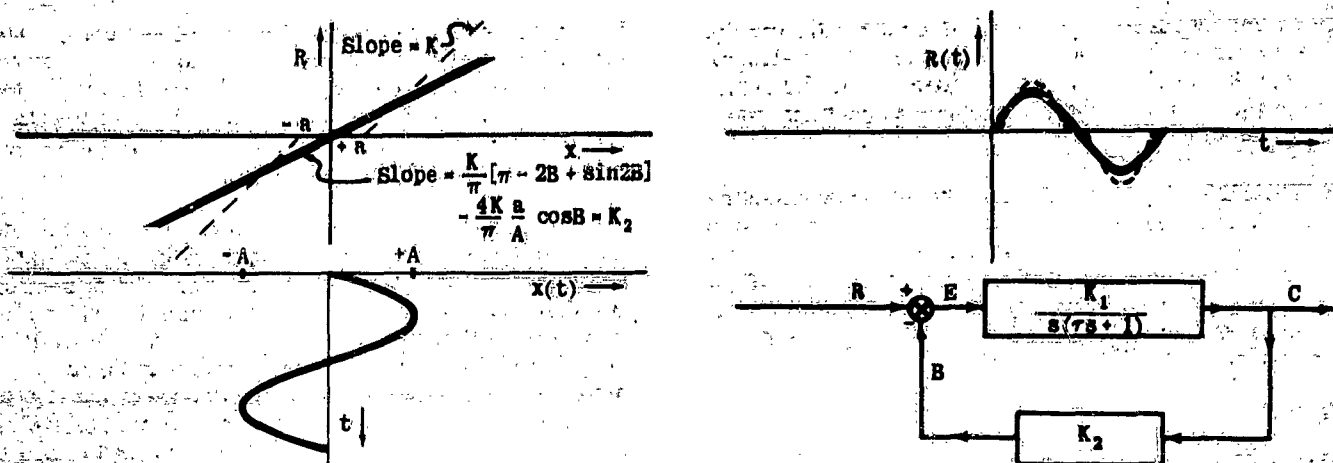


Figure VI-24. Effective Characteristics of Threshold when Harmonics are Negligible

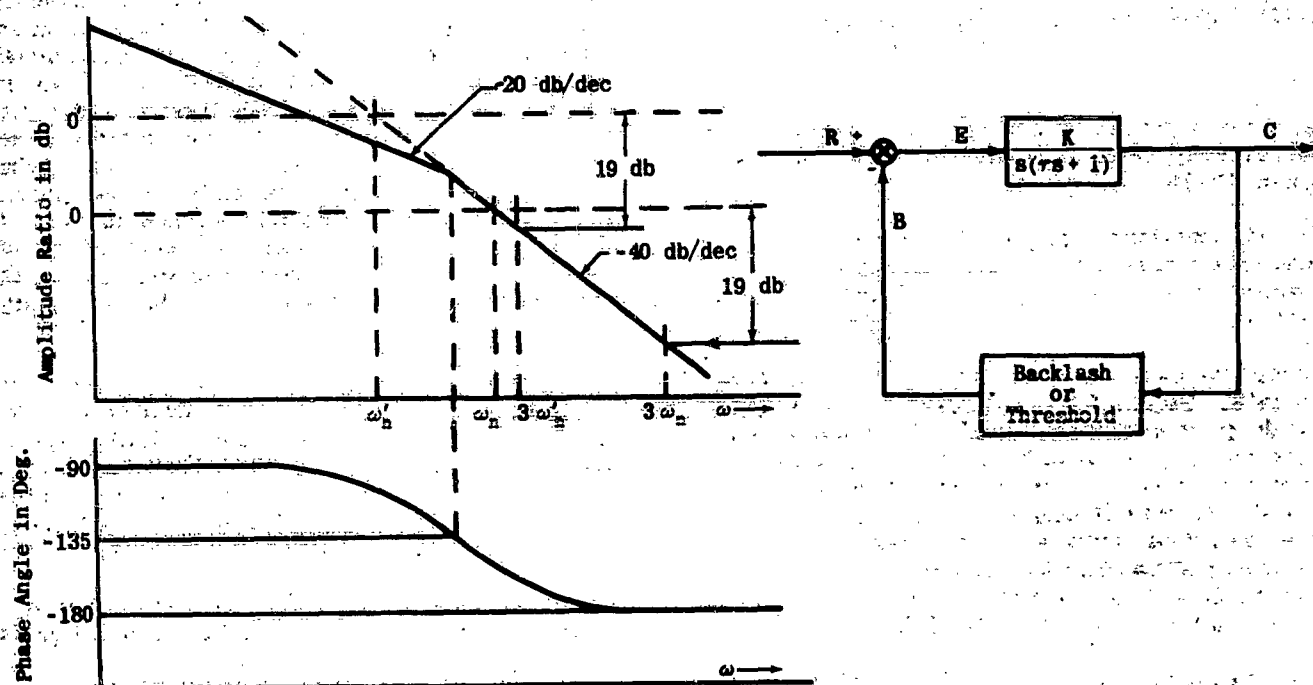


Figure VI-25. Open-Loop Bode Diagram of System for Threshold or Backlash

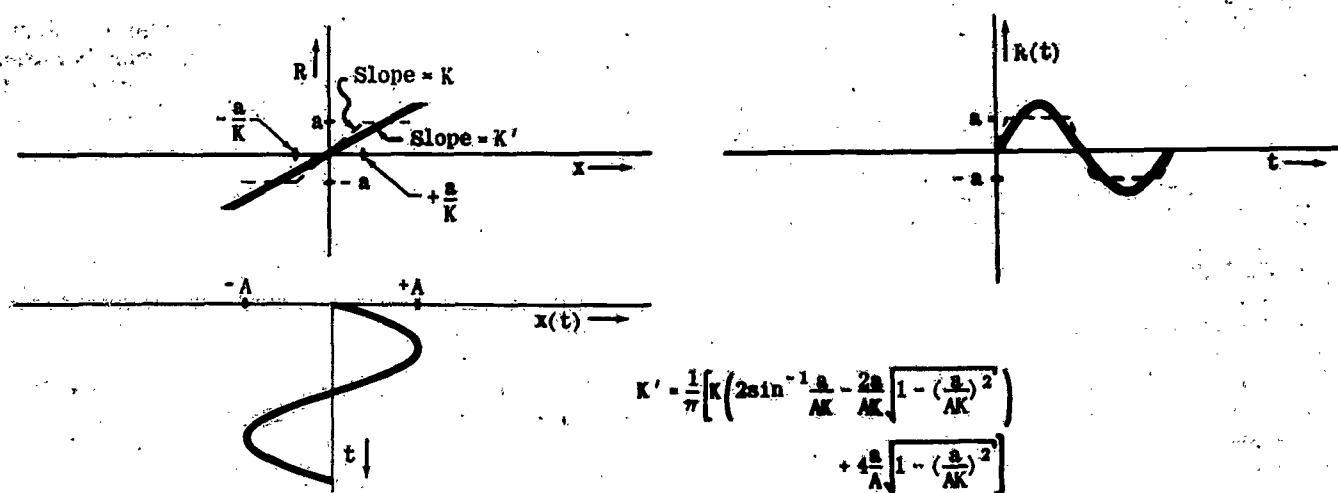


Figure VI-26. Effective Characteristics of Limiting when Harmonics are Negligible

quency and increase the apparent damping ratio for the cases in which the input amplitude to the non-linear element does not exceed the limiting value by large enough values to invalidate the assumptions concerning the harmonics.

HYSTERESIS. From the Fourier series accompanying figure VI-12, it may be seen that when the harmonics are neglected, an element containing hysteresis will have an effective shift of phase as well as a shift in amplitude. The gain and phase of the fundamental, as obtained from the Fourier expansion, are given by equations (VI-15) and (VI-16).

(VI-15)

$$\text{Gain} = \frac{1}{\pi} \left[1 - u^2 + \left(\frac{3\pi}{2} - \sin^{-1} u \right)^2 + 2 \left(\frac{3\pi}{2} - \sin^{-1} u \right) \sqrt{1 - u^2} \right] u$$

(VI-16)

$$\text{Phase} = \tan^{-1} \frac{u^2 - 1}{\left(\frac{3\pi}{2} - \sin^{-1} u + u \sqrt{1 - u^2} \right)}$$

The change in gain and phase as a function of the ratio of discontinuity amplitude to signal amplitude is plotted in figure VI-16.

Unlike the previous examples in which the effects of discontinuity were simple gain changes which could not make the system of figure VI-17 unstable, hysteresis may cause the system to become unstable.

Consider the case when backlash exists in the feedback path. Figure VI-25 illustrates the open-loop Bode diagram of the combination in which the controlled variable amplitude is much larger than the hysteresis range. In this range of amplitudes the effect of hysteresis is negligible. For smaller values of the controlled variable, figure VI-16 shows a decreased gain and a phase lag for the feedback element. From the Bode diagram it is apparent that the phase lag tends to reduce the stability of the system.

For a sufficiently small amplitude of the input to the non-linearity there will exist a combination of phase and gain change such that the zero db line cuts the Bode diagram at a frequency at which the phase passes through 180 degrees. For this particular amplitude and frequency, the system will oscillate with a constant amplitude. This type of oscillation is called a "limit cycle." If the harmonics are negligible at the limit cycle, the oscillations will be sinusoidal. For such a case, the frequency of oscillations may be obtained from the Bode diagram, and the amplitude of oscillations from the amplitude ratio curves.

(c) PHASE PLANE

In part (b) of this section, a method was discussed which permitted the inclusion of small discontinuities in the analysis and synthesis of systems. It was emphasized that if the discontinuities were not small relative to the inputs to the non-linear elements, the method could not be applied.

The investigation of systems containing large discontinuities by analytical and graphical methods is

laborious for all except the simplest of systems, since the differential equations change abruptly at the points of discontinuity. To carry out such investigations, the initial conditions for each new equation, following a discontinuity, must be determined from the previous equation at the points of discontinuity.

Since discontinuous systems of normal complexity do not lend themselves readily to analysis by the available analytical and graphical methods, this part of section two discusses a method for determining the effects of large discontinuities on simple systems. While it is recognized that few complex systems may be approximated by such simple systems, knowledge of the effects of discontinuities on these systems may sometimes provide an insight as to their behavior in the more complex systems.

When a system can be represented by second order differential equations, the state of the system at any instant following some disturbance may be completely described in terms of the dependent variable and its derivative at that instant. Since these two quantities are all that are needed, a convenient method of describing the motions of a system is to plot one against the other in a single plane called the "phase plane." The path followed by the point representing the state of the system at various instants of time then describes the complete sequence of events following the disturbance. Such a path is referred to as the "trajectory" in the phase plane.

To illustrate the techniques used with the phase plane, a simple linear system will be analyzed in detail. This will be followed by an analysis of a similar system having spring preload as a discontinuous non-linearity. The techniques used with these illustrations will then be applied to several examples in which the feedback path of a position servo is subject to discontinuities.

To introduce the phase plane, consider the undamped second order mechanical system described by equation (VI-17).

$$(VI-17) \quad m \frac{d^2 x}{dt^2} + kx = 0$$

Since time must be eliminated as a variable to obtain an equation in terms of the position and velocity alone, equation (VI-17) may be multiplied through by dx/dt and integrated with respect to time:

$$(m/2) (dx/dt)^2 + (k/2) x^2 = h$$

or

$$(VI-18) \quad \frac{(dx/dt)^2}{2h/m} + \frac{x^2}{2h/k} = 1$$

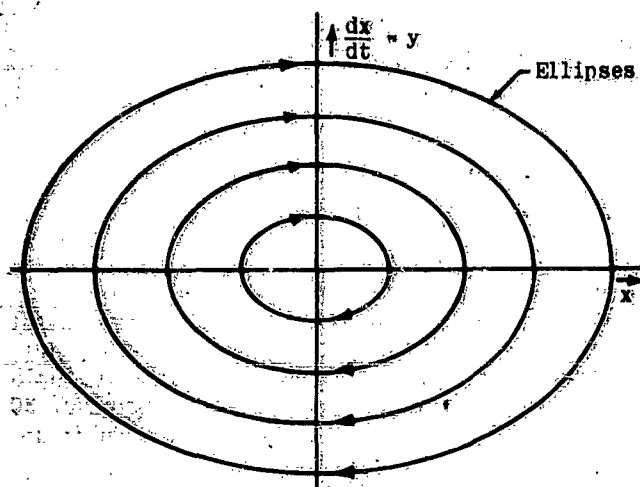
where h is the constant of integration, evaluated from the initial conditions.

If the substitutions $(dx/dt) = y$, $\alpha^2 = (2h/m)$ and $\beta^2 = (2h/k)$ are made, equation (VI-18) may be written as

$$(VI-19) \quad \frac{y^2}{\alpha^2} + \frac{x^2}{\beta^2} = 1$$

Equation (VI-19) may be recognized as that of an

ellipse with semi-axes a and b . Figure VI-27 illustrates the phase plane plot of equation (VI-19) for various values of h .



$$\text{System Equation: } m \frac{d^2 x}{dt^2} + kx = 0$$

$$\text{Equation of Trajectories: } \frac{y^2}{2h/m} + \frac{x^2}{2h/k} = 1$$

Figure VI-27. Phase Plane; Undamped Second Order System

If viscous damping is added to the mechanical system described by (VI-17), the differential equation of motion becomes (VI-20), or (VI-21).

$$(VI-20) \quad m \frac{d^2 x}{dt^2} + b \frac{dx}{dt} + kx = 0$$

$$2\zeta\omega_n = \frac{b}{m}, \quad \omega_n^2 = \frac{k}{m}$$

$$(VI-21) \quad \frac{d^2 x}{dt^2} + 2\zeta\omega_n \frac{dx}{dt} + \omega_n^2 x = 0$$

By making the substitution $(dx/dt) = y$, (VI-21) may be written in the form

$$(VI-22) \quad \frac{dy}{dt} + 2\zeta\omega_n y + \omega_n^2 x = 0$$

If equation (VI-22) is divided through by y , and rearranged, (VI-23) results,

$$(VI-23) \quad \frac{dy}{dx} = - \left(\frac{2\zeta\omega_n y + \omega_n^2 x}{y} \right)$$

in terms of the two variables x and y alone. Equation (VI-23) may be integrated to the form of (VI-24),

$$(VI-24)$$

$$y^2 + 2\zeta\omega_n xy + \omega_n^2 x^2 = Ce^{\left(\frac{2\zeta\omega_n}{\omega_n^2} \tan^{-1} \frac{y}{x} \right)}; \quad \omega_1 = \omega_n \sqrt{1 - \zeta^2}$$

but this expression is not easily used. To express the equation in a more usable form, a change of variables such as $u = \omega_1 x$ and $v = y + \zeta\omega_n x$ may be made, so that the equation may be written as (VI-25)

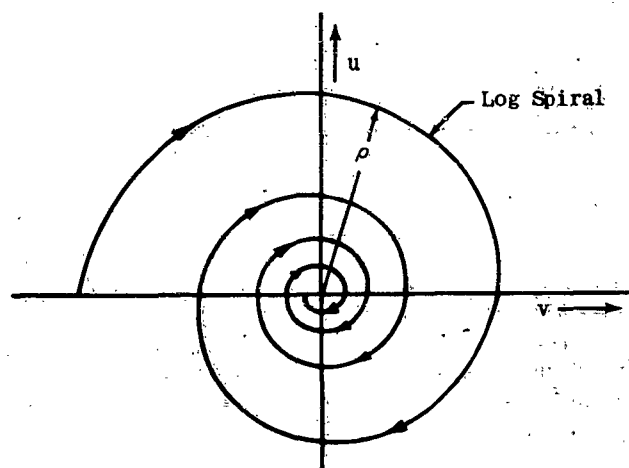
$$(VI-25) \quad v^2 + u^2 = Ce^{\left(\frac{2\zeta\omega_n}{\omega_1^2} \tan^{-1} \frac{v}{u} \right)}$$

which, in polar form, becomes (VI-26).

$$(VI-26) \quad \rho = C_1 e^{\frac{\zeta\omega_n}{\omega_1^2} \psi}$$

Equation (VI-26) is that of a logarithmic spiral in the $u-v$ plane and may be of the form illustrated in figure VI-28.

When discontinuities are introduced into the phase plane, more than one equation is needed to describe the motions of a system completely. A transformation of variables is not always possible when more than one equation is involved, without increasing the difficulty in interpreting the results. Therefore, it is often desirable to obtain the path or trajectory of the position and velocity on the phase plane using the original coordinate system. To obtain the trajectories given by equation (VI-24) without actually solving the equation, the method of isoclines can be used. In equation (VI-23) it may be noted that the slope (dy/dx) of the trajectory in the phase plane depends only upon the values of x and y . It may be further noted that in this equation the slope is constant for any particular ratio of x to y . For this particular example then, straight lines may be drawn in the phase plane corresponding to different ratios of x to y . If any trajectory intersects these lines, it must do so with the slope defined by equation (VI-23). In general, equations of the form of (VI-23) may always be used to obtain the curves through which the trajectories must pass with a particular slope. These curves are referred to as "isoclines." If a sufficient number of isoclines are drawn in a phase plane, a trajectory may be drawn by intersecting each isocline with the appropriate slope.



$$\text{System Equation: } m \frac{d^2 x}{dt^2} + b \frac{dx}{dt} + kx = 0$$

$$\text{Equation of Trajectories: } \rho = C_1 e^{\frac{\zeta\omega_n}{\omega_1^2} \psi}$$

$$\text{Where: } u = \rho \cos \psi, \quad v = \rho \sin \psi$$

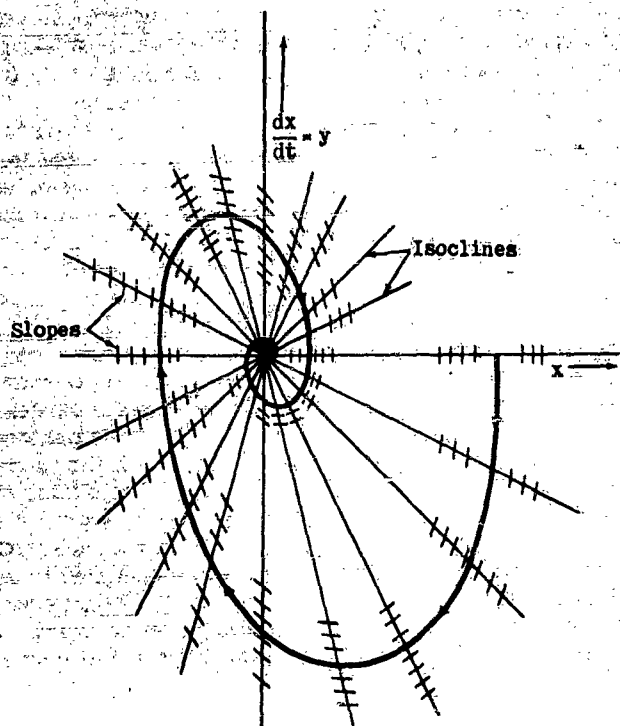
Figure VI-28. Phase Plane; Second Order System with Damping Ratio < 1

In figure VI-29 the isoclines are drawn as straight lines through the origin. The slopes, as determined from equation (VI-23), are drawn as short dashes through the isoclines to act as guides in drawing the trajectory from some initial condition.

Chapter VI

Section 3

The direction traveled by a point in a trajectory is found from the original equations. Consider, for example, the fourth quadrant of figure VI-29. In this quadrant, y is negative so that (dx/dt) is also negative; and from equation (VI-22), $(dy/dt) = -(+2\omega_n y + \omega_n^2 x)$, so that for $\omega_n^2 x > 2\omega_n y$, (dy/dt) is also negative. If the other quadrants are similarly checked, the direction of travel on the phase plane will be found as indicated.



System Equation: $m \frac{d^2x}{dt^2} + b \frac{dx}{dt} + kx = 0$
 Slope of Trajectories: $\frac{dy}{dx} = -\frac{1}{m} \left(\frac{by + kx}{y} \right)$

Figure VI-29. Phase Plane; Second Order System with Damping Ratio < 1

The simple linear systems described above serve to illustrate some of the techniques used with the phase plane representation. Now, to include a discontinuity, the spring constant term of equation (VI-17) will be changed to represent a preloaded spring. This changes the equation to the form $m(d^2x/dt^2) + (kx + f_0 \text{sgn } x) = 0$.

This relation may be expressed by the pair of equations

$$(VI-27) \quad m \frac{d^2x}{dt^2} + kx + f_0 = 0 \quad \text{for } x > 0$$

$$m \frac{d^2x}{dt^2} + kx - f_0 = 0 \quad \text{for } x < 0$$

where f_0 is the amount by which the spring is preloaded. Let $k = m\omega^2$ and $f_0 = a m \omega^2$, equations (VI-27) may be written in the form

$$(VI-28) \quad \frac{d^2x}{dt^2} + \omega^2 x + a\omega^2 = 0 \quad \text{for } x > 0$$

$$\frac{d^2x}{dt^2} + \omega^2 x - a\omega^2 = 0 \quad \text{for } x < 0$$

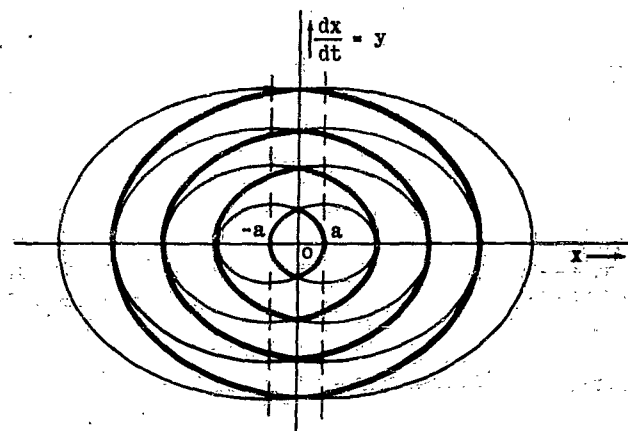
If new choices of origin are taken such that $x_1 = x + a$ for $x > 0$ and $x_2 = x - a$ for $x < 0$, equations (VI-28) may be written, after integration, as

$$(VI-29) \quad \frac{y_1^2}{2h} + \frac{x_1^2}{2h/\omega^2} = 1 \quad \text{for } x > 0$$

$$(VI-30) \quad \frac{y_2^2}{2h} + \frac{x_2^2}{2h/\omega^2} = 1 \quad \text{for } x < 0$$

where h is a constant of integration.

This system is represented in the phase plane by pairs of ellipses with origins displaced from $x = 0$ by $\pm a$ on the x axis. As a point on the trajectory passes from the left half plane into the right half plane, the trajectory changes from the ellipse at origin $x = +a$ to the ellipse at origin $x = -a$. Figure VI-30 shows the trajectories for four initial conditions.



System Equation: $m \frac{d^2x}{dt^2} + kx = -f_0 \text{sgn } x$
 $= -f_0 \text{ For } x > 0$
 $= +f_0 \text{ For } x < 0$

Equations of Trajectories: $\frac{y_1^2}{2h} + \frac{x_1^2}{2h/\omega^2} = 1 \text{ For } x > 0$
 $\frac{y_2^2}{2h} + \frac{x_2^2}{2h/\omega^2} = 1 \text{ For } x < 0$

Where: $x_1 = x + a$; $x_2 = x - a$

Figure VI-30. Phase Plane; Second Order with Spring Preload

As illustrated by the simple undamped system with a preloaded spring, the inclusion of discontinuities in the phase plane representation of systems may be simply achieved. The labor of determining the initial conditions at the points of discontinuity is eliminated, since the velocity and position are always available for all points in the plane including points of discontinuity.

The remainder of this section will present several examples of simple systems containing discontinuities to further illustrate the phase plane method.

COULOMB FRICTION. The effects of coulomb friction

upon a mass-spring combination will now be considered. The equation relating the mass-spring friction combination is that of (VI-31).

$$(VI-31) \quad m \frac{d^2x}{dt^2} + kx = -F \operatorname{sgn} \frac{dx}{dt} = F'$$

If the substitutions $F' = -f_0$ for $(dx/dt) < 0$, $F' = f_0$ for $(dx/dt) > 0$, $k = m\omega^2$ and $f_0 = am\omega^2$ are made, (VI-31), may be written in the form of equations (VI-32) and (VI-33), which, upon integrating, become those of (VI-34) and (VI-35)

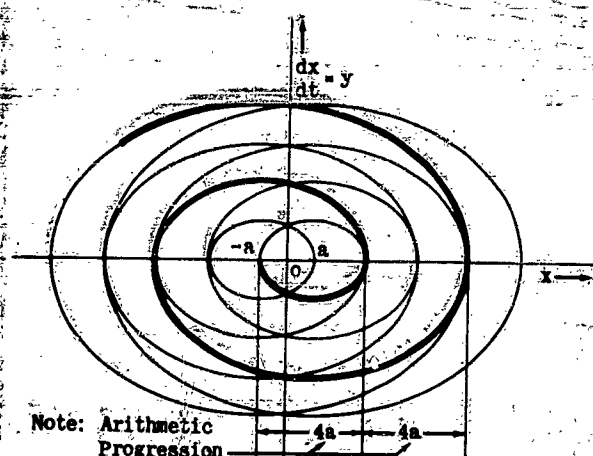
$$(VI-32) \quad \frac{d^2x}{dt^2} + \omega^2(x - a) = 0 \quad \text{for } \frac{dx}{dt} < 0$$

$$(VI-33) \quad \frac{d^2x}{dt^2} + \omega^2(x + a) = 0 \quad \text{for } \frac{dx}{dt} > 0$$

$$(VI-34) \quad \frac{y_1^2}{2h} + \frac{x_1^2}{2h/\omega^2} = 1 \quad \text{for } \frac{dx}{dt} < 0$$

$$(VI-35) \quad \frac{y_2^2}{2h} + \frac{x_2^2}{2h/\omega^2} = 1 \quad \text{for } \frac{dx}{dt} > 0$$

These latter two equations are those of ellipses in the phase plane.



System Equation: $m \frac{d^2x}{dt^2} + kx = -F \operatorname{sgn} \frac{dx}{dt}$

$= +f_0$ For $\frac{dx}{dt} < 0$

$= -f_0$ For $\frac{dx}{dt} > 0$

Equations of Trajectories:

$$\frac{y_1^2}{2h} + \frac{x_1^2}{2h/\omega^2} = 1 \quad \text{For } \frac{dx}{dt} < 0$$

$$\frac{y_2^2}{2h} + \frac{x_2^2}{2h/\omega^2} = 1 \quad \text{For } \frac{dx}{dt} > 0$$

Figure VI-31. Phase Plane; Second Order System with Coulomb Friction

In this case the trajectories follow the ellipses with center at $x = -a$ for positive values of (dx/dt) and the

ellipses with center at $x = +a$ for negative values of (dx/dt) , as shown in figure VI-31. It is of interest to note that the maxima form an arithmetic progression with difference $-4a$ in contrast to the geometric progression of linear systems. Also, it is evident in figure VI-31 that a steady state value other than zero is possible.* If the methods used with small discontinuities were applied to this problem, the possibility of a steady state value other than zero would not be apparent, since that method approximates coulomb friction with an effective variable viscous friction.

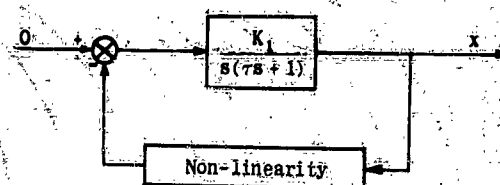
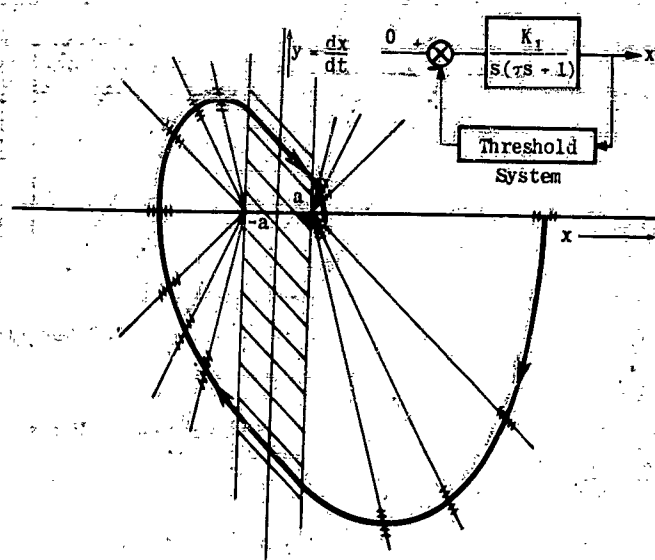


Figure VI-32. Closed-Loop System

To illustrate the effects of discontinuous non-linearities on closed-loop systems, the positional servomechanism illustrated in figure VI-32 will be discussed for three types of discontinuities in the feedback path. It will be assumed in each case that the input is zero and that the motion of the output following some initial condition is to be determined.



System Equations: Slope of Trajectories:

- (1) $\frac{d^2x}{dt^2} + A \frac{dx}{dt} = 0$ for $-a < x < a$ $\frac{dy}{dx} = -A$ for $-a < x < a$
- (2) $\frac{d^2x_1}{dt^2} + A \frac{dx_1}{dt} + Bx_1 = 0$ for $x > a$ $\frac{dy}{dx} = -\left(\frac{Ay_1 + Bx_1}{y_1}\right)$ for $x > a$
- (3) $\frac{d^2x_2}{dt^2} + A \frac{dx_2}{dt} + Bx_2 = 0$ for $x < -a$ $\frac{dy}{dx} = -\left(\frac{Ay_2 + Bx_2}{y_2}\right)$ for $x < -a$

Figure VI-33. Phase Plane; Closed-Loop System with Threshold in Feedback

* At this value the spring restoring force is not greater than the coulomb friction force, and no further motion is possible.

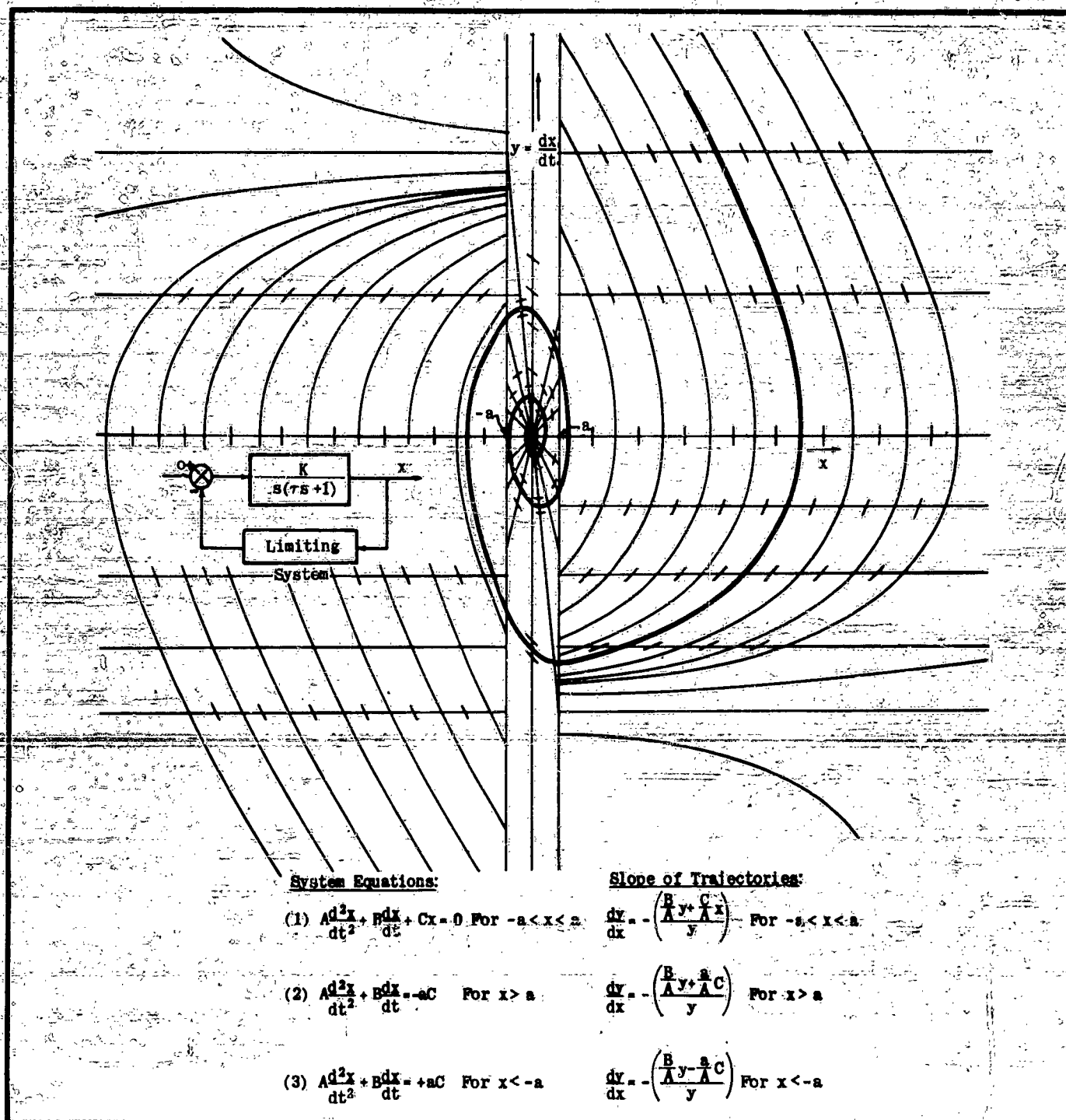


Figure VI-34. Phase Plane; Closed-Loop System with Limiting in Feedback

THRESHOLD. When the output is fed back through a component having threshold, it must exceed the threshold limits a and $-a$ (see figure VI-5) before any signal reaches the input. To describe this system for all values of the output quantity, equations (VI-36), (VI-37), and (VI-38) can be written.

$$(VI-36) \quad \tau \frac{d^2x}{dt^2} + \frac{dx}{dt} = 0 \quad \text{for } -a < x < a$$

$$(VI-37) \quad \tau \frac{d^2x}{dt^2} + \frac{dx}{dt} + k_1 k_2 (x - a) = 0 \quad \text{for } x > a$$

$$(VI-38) \quad \tau \frac{d^2x}{dt^2} + \frac{dx}{dt} + k_1 k_2 (x + a) = 0 \quad \text{for } x < -a$$

where x is the quantity fed back and k_2 is the slope of the threshold curve. By making the substitutions $x_1 = x - a$ for $x > a$ and $x_2 = x + a$ for $x < -a$, these equations may be written, after rearranging, in the form of (VI-39), (VI-40), and (VI-41).

$$(VI-39) \quad \frac{d^2x}{dt^2} + A \frac{dx}{dt} = 0 \quad \text{for } -a < x < a$$

$$(VI-40) \quad \frac{d^2x_1}{dt^2} + A \frac{dx_1}{dt} + Bx_1 = 0 \quad \text{for } x > a$$

$$(VI-41) \quad \frac{d^2x_2}{dt^2} + A \frac{dx_2}{dt} + Bx_2 = 0 \quad \text{for } x < -a$$

where $A = 1/\tau$ and $B = (k_1 k_2)/\tau$. If the substitution $y = (dx/dt)$ is made in (VI-39) and this relation divided through by y , it may be written in the form of (VI-42).

$$(VI-42) \quad \frac{dy}{dx} = -A \quad \text{for } -a < x < a$$

That is, for values of x between a and $-a$, the trajectories in the phase plane will have a constant slope of magnitude $-A$. Equations (VI-40) and (VI-41) are of the same form as equation (VI-21) and result in equations similar to equation (VI-23), or

$$(VI-43) \quad \frac{dy_1}{dx_1} = -\left(\frac{Ay_1 + Bx_1}{y_1}\right) \quad \text{for } x > a$$

$$(VI-44) \quad \frac{dy_2}{dx_2} = -\left(\frac{Ay_2 + Bx_2}{y_2}\right) \quad \text{for } x < -a$$

The latter three equations provide the necessary information for the drawing of the isoclines in the phase plane.

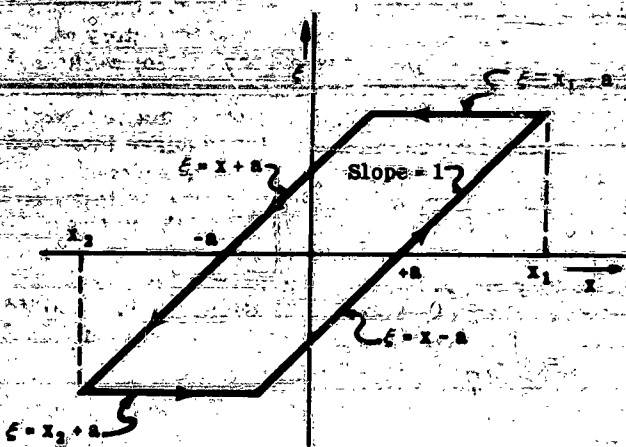


Figure VI-35. Hysteresis Curve

The trajectory shown in figure VI-33 was drawn with the aid of these isoclines. As indicated by this figure, the presence of threshold in the feedback path of a position servomechanism may cause a steady state error equal in magnitude to the threshold value. When the method used with small discontinuities was applied to this problem, the possibility of a steady state value other than zero was not apparent, since that method approximated the non-linear transfer characteristic by a straight line.

LIMITING. When the quantity fed back in figure VI-32 is limited so that values of the output exceeding $\pm a$ will not be fed back, three equations (VI-45), (VI-46), and (VI-47) are again needed to describe all the possible motions of the system.

$$(VI-45) \quad A \frac{d^2x}{dt^2} + B \frac{dx}{dt} + Cx = 0 \quad \text{for } -a < x < a$$

$$(VI-46) \quad A \frac{d^2x}{dt^2} + B \frac{dx}{dt} = -aC \quad \text{for } x > a$$

$$(VI-47) \quad A \frac{d^2x}{dt^2} + B \frac{dx}{dt} = aC \quad \text{for } x < -a$$

The corresponding equations for the isoclines are:

$$(VI-48) \quad \frac{dy}{dx} = -\left(\frac{By + Cx}{y}\right) \quad \text{for } -a < x < a$$

$$(VI-49) \quad \frac{dy}{dx} = -\left(\frac{By + aC}{y}\right) \quad \text{for } x > a$$

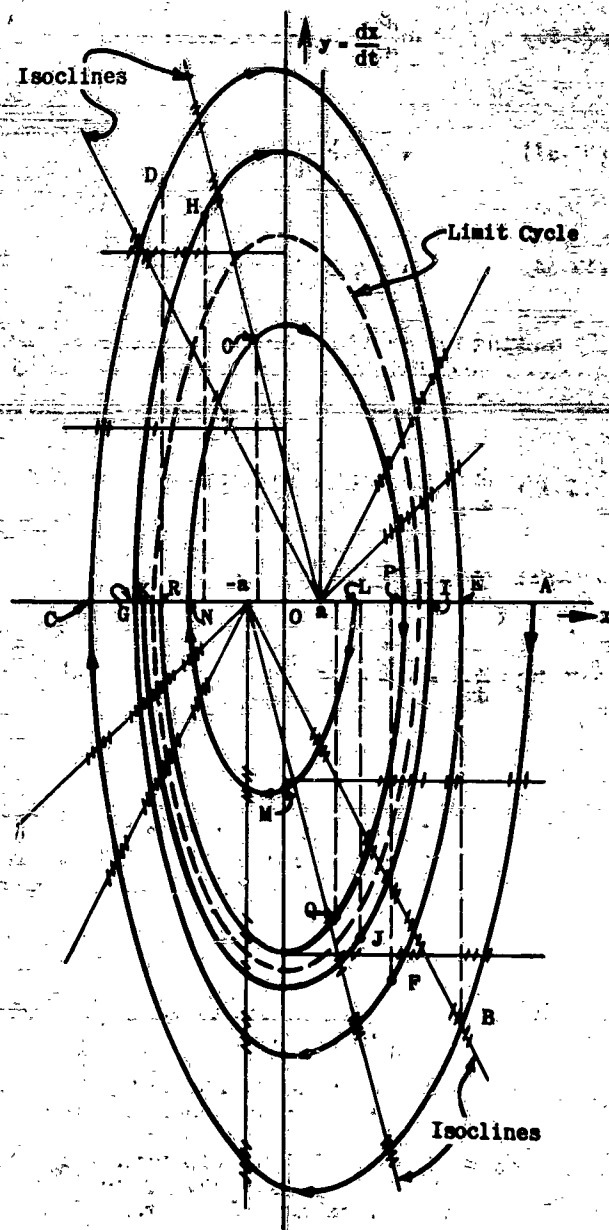


Figure VI-36. Phase Plane; Closed-Loop System with Hysteresis in Feedback

$$(VI-50) \quad \frac{dy}{dx} = -\left(\frac{B}{A}y - \frac{C}{A}\right) \quad \text{for } x < -a$$

These three equations with $A = B = 1$ and $C = 10$ were used to draw the isoclines of figure VI-34.

HYSTERESIS. Figure VI-35 illustrates the hysteresis curve which describes the relation between the output (x) and the quantity fed back (z); x_1 and x_2 are the maximum amplitudes reached by the output, and $2a$ is the maximum range through which the output can move without causing any change in the quantity fed back.

The differential equations for the system are

$$(VI-51) \quad A \frac{d^2x}{dt^2} + B \frac{dx}{dt} + C(x - a) = 0$$

$$(VI-52) \quad A \frac{d^2x}{dt^2} + B \frac{dx}{dt} + C(x_n - a) = 0$$

$$(VI-53) \quad A \frac{d^2x}{dt^2} + B \frac{dx}{dt} + C(x + a) = 0$$

$$(VI-54) \quad A \frac{d^2x}{dt^2} + B \frac{dx}{dt} + C(x_n + a) = 0$$

By making the substitutions $y = (dx/dt)$, $x_1 = x - a$, $x_2 = x + a$ and dividing the resulting equations by y , these become

$$(VI-55) \quad \frac{dy_1}{dx_1} = -\left(\frac{B}{A}y_1 + \frac{C}{A}x_1\right) \quad \text{with center at } x = a$$

$$(VI-56) \quad \frac{dy}{dx} = -\left(\frac{B}{A}y + \frac{C}{A}D_1\right) \quad \text{where } D_1 = (x_n - a)$$

$$(VI-57) \quad \frac{dy_2}{dx_2} = -\left(\frac{B}{A}y_2 + \frac{C}{A}x_2\right) \quad \text{with center at } x = -a$$

$$(VI-58) \quad \frac{dy}{dx} = -\left(\frac{B}{A}y + \frac{C}{A}D_2\right) \quad \text{where } D_2 = (x_n + a)$$

Equations (VI-55) and (VI-57) are of the same form as previously discussed. Equations (VI-56) and (VI-58), however, are functions of x_n . Therefore, no family of isoclines can be drawn until an initial value of x_n is obtained. For the purpose of drawing the trajectories in figure VI-36, the coefficients of the equations and the hysteresis range were assigned the values $A = B = a = 1$, and $C = 10$.

To start the construction of the trajectory, an initial value of x_n was chosen at point A in the figure. With this value of x_n , equation (VI-56) provided the slopes of the trajectory for different values of y , which permitted drawing the trajectory from A to B. Point B is a distance of $2a$ (total hysteresis range) on the x -axis from A, and is the point at which the equations of slopes must change from (VI-56) to (VI-57). From point B to C the trajectory follows the slopes given by equation (VI-57). The value of x at point C is the x_n of equation (VI-58) which is valid until the trajectory reaches point D. From point D to E equation (VI-55) provides the slopes of the trajectory. At point E the process is repeated except for the different values of x_n used as the trajectory passes through the x -axis.

In figure VI-36 a second starting point was taken at point L and the above procedure repeated. Unlike the other closed-loop examples, the steady state motion of the output is other than zero. The trajectory starting at A tends to spiral toward the origin; while that starting point nearer the origin, point L, tends to spiral away from the origin. Both of these spirals are actually approaching the closed curve indicated by a dashed line. The significance of this closed curve is that the system will operate in what is termed a limit cycle. That is, irrespective of the initial conditions the steady state output of the system will be oscillatory and of a magnitude indicated by the dashed line. That steady state oscillations could occur in such a system has been shown by the methods used with small discontinuities.

BIBLIOGRAPHY

The following bibliography is included for reference. The list is in no sense complete, but contains the major source material for this chapter. Many of the references, themselves, contain much more complete and detailed bibliographies.

1. 'Problème général de la stabilité du mouvement,' by M. A. Liapounoff; Annals of Mathematical Studies, Vol. 17, Princeton Press.
2. 'Introduction to Non-Linear Mechanics,' by N. Minorsky; J. W. Edwards, Ann Arbor, 1947.
3. 'Small Discontinuous Non-Linearities,' by D. T. McRuer and R. G. Halliday; Northrop Aircraft, Inc., Unpublished paper, 1952.
4. 'The Effects of Backlash and of Speed-Dependent Friction on the Stability of Closed-Cycle Control Systems,' by A. Tustin; The Journal of the Institution of Electrical Engineers, Vol. 94, Part II A, No. 1, May, 1947.
5. 'The Dynamics of Automatic Controls,' by R. C. Oldenbourg and H. Sartorius; The American Society of Mechanical Engineers, New York, 1948.

6. 'Analysis of Relay Servomechanisms,' by H. K. Weiss; Journal of the Aeronautical Sciences, Vol. 13, No. 7, July 1946.
7. 'The Analysis of Relay Servomechanisms,' by D. A. Kahn; Report No. R 48-22., Curtiss-Wright Corp., Columbus, Ohio, December, 1948.
8. 'Theory of Oscillations,' by A. A. Andronow and C. E. Chaikin; Princeton University Press, New Jersey, 1949.

CHAPTER VII

MACHINE METHODS

SECTION 1 - INTRODUCTION

In the preceding chapters of this volume various methods of analysis and synthesis have been discussed. It has been pointed out that the matter of primary interest in control systems work is the transient response; and the methods previously used have been considered from the point of view of obtaining approximations to that response.

All of these methods have certain defects, but are used to obtain a sufficiently close approximation to the transient. What is meant by sufficiently close at any point of the design process is a matter of judgment, but the essential point is that as design continues and more and more exact decisions have to be made, more and more complete and accurate information about the transient is needed.

The methods heretofore discussed have been arranged in a hierarchy of utility. These methods are desirable not only from the standpoint of usefulness, but also because they help the designer develop a "feel" for the physical situation. These methods, however, leave certain things to be desired. For example, if the system under consideration is a multi-loop affair, it becomes difficult and time-consuming to isolate the effects of varying an inner-loop parameter on the overall behavior of the entire system. Systems of this type are frequently used to control the behavior of aircraft, which are, in themselves, dynamical systems of no mean order of complexity.

In order to handle complex multi-loop systems some method is needed which rapidly determines the total transient response in such a way that the effects of varying any parameter are easily isolated and observed. Theoretically, it would be possible, say by hand computation of enough individual cases, to find an optimum combination of parameters. However, from the point of view of the time required to complete a design, this procedure is often outside practical consideration. The result is that a thorough analysis is not performed, and an extensive period of debugging the actual physical system is required. For this

reason, it now becomes necessary to consider what can be done by automatic means of computation. This will be the general subject matter of this chapter.

In the following section, the reasons for using machines in analysis will be discussed more thoroughly than has been done in the above cursory and purely introductory examination. The potentialities and limitations of these devices will be examined. Special attention will be given to problems which, to date, can be solved most practicably with computer techniques.

It will be shown that there are only two types of devices suitable for such purposes at present; the analog, and the digital type of computer. The two types of computer (and variations within each type) will be discussed, primarily with a view to indicate which are best for control systems analysis and synthesis, and why they are best.

During this discussion, it will become apparent that the operational amplifier type of analog computer appears to be most useful for the applications considered here. The final section of this chapter will attempt to establish that this proposition is indeed true, and to say under what conditions, and for what type of problems this superiority of the operational amplifier analog most manifests itself.

In order to do this, it will be necessary to describe, in some detail, how the computer works. This section will therefore lay the groundwork for the following chapter, in which the properties and operation of the operational amplifier computer will be considered in detail. Section 3 of this chapter will thus be in some sense preliminary, but enough material will be given to form a firm basis for the discussion of the relative advantages of this type of computer over other kinds. Special attention will be given to such qualities as easy and rapid variation of basic parameters, ability to collect data in a form suited for immediate inspection, and usefulness in simulation while using actual system components.

SECTION 2 - NEED FOR MACHINE METHODS

The dynamical problems presented by the design of aircraft control systems are complex. They involve many degrees of freedom, and may contain a multiplicity of feedback loops. Some of these feedbacks are inherent in the dynamics of the airframe itself; others are

added, in the design of the control system.

Multiple-loop systems give rise to the most difficult problems in control systems design, primarily because of the difficulty of determining how a change in an inner

Chapter VII

Section 3

loop parameter affects the overall transient response of the entire system. An attempt is sometimes made to avoid the problem of determining these effects. For example, one part of the system may be overdesigned and forced to meet too stringent requirements in an effort to assure that unknown requirements on inner loop parameters are met.

The techniques previously discussed do not adequately meet the problem discussed above. The analytical methods replace a complicated system by a single closed loop transfer function for the complete subsystem. This closed loop expression often contains the parameters of the components in complicated combinations rather than as individual and easily isolated quantities. When such expressions are then used as part of a dynamical relation, it becomes extremely difficult to determine the effects due to any single parameter of a component in the inner loop.

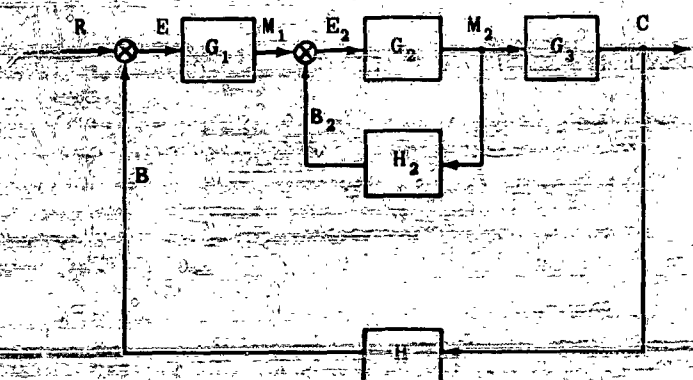


Figure VII-1. Illustrative Two-Loop Servomechanism

For example consider the system shown in figure VII-1; it is a fairly simple two-loop device. Its transfer function is:

$$(VII-1) \quad \frac{C}{R} = \frac{G_1 G_2 G_3}{1 + G_2 H_2 + G_1 G_2 G_3 H}$$

It is evident from the form of this expression that even for this comparatively simple case it would not be easy to say anything about the effects of the individual parameters contained in G_1 , or H_2 , upon the behavior of the entire system.* In more involved multiple-loop systems, such as found in aircraft control systems, the determination of effects of inner loop parameter variation becomes an horrendous task and is frequently economically unfeasible. The technique of overdesign (which was discussed above) may be applied, and the efficient design of the entire controls system, as an integrated unit, is imperilled.

To resolve this difficulty, some means of analysis is needed which will permit the effects of varying all parameters to be quickly and easily observed, no matter how the component bearing these parameters is interconnected with the rest of the system.

It is not to be understood that the reasons discussed above are the only ones which indicate that it would be desirable to have a means by which parameter changes could be readily handled. The analysis of any complex system, with a number of adjustable parameters, would be greatly expedited by some such method. Even after a large amount of preliminary thinking has been done about suitable (and attainable) ranges of adjustable quantities, there may well remain a large amount of analysis to be done before the most suitable combination of parameters can be selected.

For these reasons it now becomes necessary to investigate machine methods. The next section will discuss the various types of computer and the basic mode of operation for each. The final section will compare the types of computers and their relative usefulness in solving problems of the sort discussed in this section.

SECTION 3 -- AVAILABLE METHODS OF AUTOMATIC COMPUTATION

Section 2 of this chapter has shown that a need exists for some automatic means of solving the equations of motion of complex systems. This section will consider two different ways in which this may be done by the use of automatic computation. The two types of computation considered are the analog and the digital methods.

(a) DIGITAL COMPUTERS

A digital computer is any device which solves mathematical problems by the numerical process of counting discrete quantities. The digital computer may be a device as simple as the ancient abacus or as complicated as the modern electronic giant "Whirlwind I," but the fundamental principles of operation are the same.

The normal desk computing machine is another example of a digital computer. Since this is a relatively simple device, it will be used to illustrate, in somewhat more concrete form, the general principle of operation stated above.

Suppose that it is required to draw a graph of the equation $y = 5x + 10$. An operator would choose discrete values of x , say x_1, x_2, \dots, x_n and use the machine to multiply each x_i by 5. The machine actually operates in the following way:

The quantity one is added x_1 times, and then this sum is added five more times. In other words the machine has operated by counting the discrete quantity one, $5x_1$ times. Another interesting property can be noted: The machine stored the sum of x_1 units so that the second operation of addition could take place later.

After each x_i was multiplied by 5, the operator would

* To do so, of course, requires at least a knowledge of their effects upon the roots of the characteristic equation. Many of the methods previously developed are essentially ways of getting this information without actually factoring that equation, but it is evident from (VII-1) that these methods will not work here, because of the way in which G_2 , H_2 (and hence their parameters) are involved.

write down the resulting product, and at some later time add 10 to obtain y_1 . In this operation, the operator stored the information $5x_1$ so that the quantity 10 might be added at a later time. When all the y_i were computed and tabulated, the graph $y = f(x)$ could be drawn. The precision of the graph would depend on the maximum number of decimal places available in the machine, since this would determine the minimum possible interval between the x_i . Obviously the graph could be made as accurate as desired, within the limits imposed by the available number of decimal places, by choosing sufficiently small intervals between the x_i .

The complete procedure and operation necessary to solve the problem could be outlined in the following way:

1. Operator decided on some plan of action or "program."
2. Operator fed some information into the machine.
3. Machine operated on some information, and stored some information.
4. Operator stored some information and fed some new information into the machine.
5. Machine operated on new information and produced answers that could be tabulated.
6. Operator drew graph from tabulated information.

Equations of the types considered in this volume could always be solved numerically if time were available; solutions to these equations might be obtained by using a desk computing machine, somewhat in the manner described above. This would, of course, be very time consuming. It is to overcome this objection that modern high speed digital computers have been built. In these machines, both numerical information and programming are introduced at the beginning of a problem. All required operations (including storage) then become completely automatic. Final answers generally are received in the form of tabulations or punched cards, but new methods have been developed to plot answers in graphical form also.

Many techniques have been developed to speed up computation. Electronics has been used to replace mechanical parts. Counting systems other than decimal are frequently used. These details will not be discussed here, but are merely mentioned as matters of background interest.

(b) ANALOG COMPUTERS

An analog computer is a physical system the variables of which may be easily measured, controlled or manipulated; it is used to study another physical system which does not have these desirable characteristics. In addition, the physical system represented by the analog computer must be governed by the same mathematical relationships as the system under study.

The analog computer may be a device as simple as a slide rule in which a length along the rule is made proportional to the logarithm of the number concerned. On the other hand the analog computer may be as complex as the modern electronic computer whose electrical system is set up according to the mathematical relation-

ships of the system to be studied.

There are many types of analog computers, but they all may be divided into two categories:

1. Those devices in which the detailed structure of each component of the system under investigation is represented as its analog (e.g., a spring is made analogous to a condenser, a damper to a resistance).
2. Those devices which function by performing the mathematical operations indicated by the differential equations representing the systems to be studied.

An example of the first type is a transient analyzer. This device might be set up to investigate the dynamics of a mechanical system which contains a mass, M , damping, B , and a spring constant, k , and is described by the differential equation:

$$(VII-2) \quad M \frac{d^2x}{dt^2} + B \frac{dx}{dt} + kx = f(t)$$

One set-up which might be used as analogous to this mechanical system is an electrical L, R, C series circuit excited by a voltage $e(t)$ and which has as its equation

$$(VII-3) \quad L \frac{d^2q}{dt^2} + R \frac{dq}{dt} + \frac{1}{C} q = e(t)$$

The mathematical form of the expressions in (VII-2) and (VII-3) is identical provided that the inputs $f(t)$ and $e(t)$ depend upon the time, t , in the same way functionally. Thus the behavior of the electrical system can be studied, and the results applied to the performance of the mechanical system.

The magnitude of the inductance may be made numerically equal to the magnitude of the mass; the resistance, to the damping; and the inverse of the capacitance, to the spring constant.

The second type of analog computer is one in which devices are used essentially to perform mathematical operations as such, rather than to mimic more directly the behavior of the system being represented.

These computers are of the "differential analyzer" electromechanical type, using, for example, Kelvin ball-and-disk integrators, and of the "operational amplifier" type. These operational amplifiers perform the mathematical operations of addition, multiplication by constants, and integration with respect to time. In addition, special apparatus may be used for multiplication of variables and the introduction of discontinuity type non-linearities.

In this type of computer, variables are represented by voltages; parameters are adjusted by plug-in resistors, potentiometers, and capacitors, and results are easily recorded in graphic form by means of oscillographs.

The previous example may be used to show how such a computer might be set up. Equation (VII-2) may be rewritten as:

$$M \frac{d^2x}{dt^2} = f(t) - B \frac{dx}{dt} - kx$$

The computer might be set up as indicated in the sche-

matic of figure VII-2.

Note that the only operations used are addition, multiplication by constants (including -1) or integration with respect to time. A voltage equal to $f(t)$ is fed into the computer and x and its first two derivatives may be picked off as voltages or applied to a recorder for graphical representation.

In this example the computer is easily designed to operate in the same time scale as the physical system represented. The operational amplifier type analog computer can almost always be used on this one-to-one time relationship with physical systems such as servo-mechanisms. For the special cases where frequencies involved are too high or too low for satisfactory operation, time scale changes may be easily made.

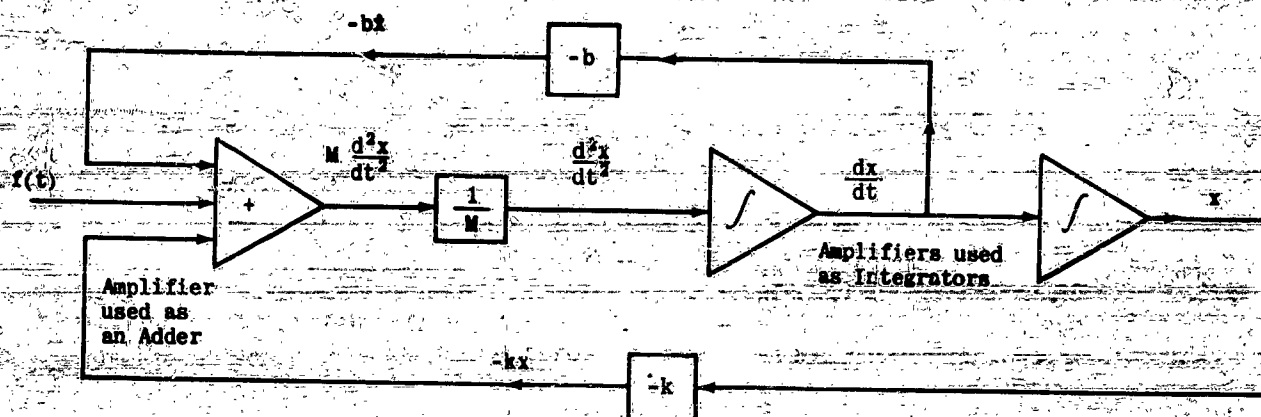


Figure VII-2. Illustrative Analog Computer Setup

SECTION 4 - EVALUATION OF MACHINE METHODS FOR CONTROL SYSTEM WORK

The previous sections of this chapter have shown the need for machine methods and the types of computers available. This section will discuss the relative merits of analog and digital computers when used in control system work. The two types will be compared and contrasted so that some conclusion may be drawn regarding the choice of computer type for specific applications.

The choice of type of computer will be dependent largely on the following items:

1. Set-up time.
2. Easy and rapid variation of parameters.
3. Ability to collect data in a form suited for immediate inspection.
4. Usefulness in simulation while connected to actual system components.
5. Accuracy.

(a) SET-UP TIME

In the last section, it was shown that digital computer operation depended upon previous programming. In complex problems, of the type treated in aircraft control work, the preliminary programming may become a Herculean task. Recall the operations required for the simple example given in the last section, and then consider the large programming necessary to solve a simple sixth order algebraic equation. This task is multiplied many times when one considers operations such as integration, differentiation and inputs such as sine waves, or irregular waves.

Many digital computers are built as "general purpose" machines. That is, they can be set-up to perform a wide variety of combinations of elementary operations necessary for the solution of various kinds of problems.

These "general" computers have the advantage of being applicable to a very large number of types of problems, but a considerable price has to be paid for this flexibility. Because of it, extensive programming and a large set-up time becomes necessary to prepare the machine for the solution of specific problems. Digital computers have also been built for special purposes but not all exigencies can be adequately covered by such "special purpose" machines, and even with "special purpose" machines, careful programming may be required.

For those types of problems where the same cycle of operations is performed repetitively, the set-up time for a digital machine becomes a smaller proportion of the total time allotted for a task. Use of analog computers may also require careful planning. Certain precautions, important in any electrical circuitry, must be observed. Details such as impedance matching, voltage levels, and efficient use of amplifiers must be examined. These considerations are of considerable more importance in the physical structure type of analog computer than in the operational amplifier type.

With both analog and digital type computers, considerable interconnecting must be done. The interconnecting must be checked and this takes time. However, the wiring required in the analog computer is considerably less than that used in the digital type because of complex programming required by the latter.

(b) VARIATION OF PARAMETERS

It is not usually easy to change the numerical values of parameters in digital computers. These machines may perform individual operations at extremely high speed. However, an entire computation may have to be

carried through one or two hundred times before the equivalent of one second time of operation of the device represented by equations is complete. In addition the results are usually presented in tabular or punch card form. For these two reasons - length of time and poor form of output data - it is difficult for an operator to evaluate quickly a result so that he may sensibly vary a parameter.

In the physical structure type analog computer, each parameter is usually represented by easily identifiable individual or group resistors, capacitors, or inductors. These components may be varied within limits, but the problem of size (mentioned earlier) must be reckoned with.

Parameters may usually be varied on the operational amplifier type by merely changing a plug-in resistor or the setting of a potentiometer. It will be shown later that all electrical analog computers present data in a form that makes evaluation simple and rapid. This is of great assistance to the operator who wants to make reasonable decisions when varying parameters.

(c) FORM OF DATA

The results of digital computers are usually presented in the form of punched cards or tabulation of numbers. As has been pointed out previously, it is difficult to visualize the behavior of a physical system from tables or punched cards. There are some new machines which convert the results into graphical form. However, even with these machines it takes a relatively long period of time to plot the equivalent of one second operation of a physical system.

The results obtained from analog computers are presented in the form of voltages. It is a simple matter to connect any output of the computer to one of several types of recorders which automatically plot a variable against time. This, together with the fact that the analog computers can work in real time, makes the evaluation of the system under investigation simple and rapid.

(d) SIMULATION

It is frequently desirable to simulate operation of a complete system by combining actual physical subsystems with the mathematical representation (computer operation) of the balance of the system. In this way, it becomes possible to avoid errors which might result from an attempt to represent the subsystem in purely mathematical form. In some cases the physical laws which govern the behavior of that subsystem have not been well worked out in mathematical form, or they are sufficiently complicated that reduction to suitable machine form might demand excessive time.

Digital computers do not operate on a real time scale, in general, and also generally do not put out signals in a form directly usable by the physical equipment;

therefore, it normally cannot be used for the purposes of simulation. The analog computer, and in particular the operational amplifier type, does work in real time and is admirably suited to simulation studies.

(e) ACCURACY

Theoretically, the results of digital computing may be made as precise as desired. The precision is limited only by the number of digits which the machine can hold in the various registers of its arithmetic unit and its memory or storage function. This limitation is more severe than is apparent at first glance. Numerical errors caused by rounding off figures and truncating numerical series tend to cumulate statistically. For this reason, it may be necessary, for example, to carry a large number of significant figures in a computation in order to be assured that the result is correct to a relatively few significant figures. This sort of difficulty may be minimized by very careful planning and programming.

The precision of analog computers is largely determined by the precision of the components within the computer. With careful planning, most analog computers can produce results accurate to within a few percent.

It should be noted that in controls system work, the parameters of components of the control system are seldom known to better than about ten percent. In the case of certain aerodynamic quantities, uncertainties as much as $\pm 50\%$ are not uncommon. For these purposes extreme accuracy is not usually justified in automatic machines.

For those special cases where an investigation requires extreme accuracy only digital machines can be used at the present time.

In conclusion, it may be stated, on the basis of the above considerations, that the operational amplifier type analog computer is generally the "best" machine for use in control systems work. In special cases, the other types should be considered.

It may not be inapposite here to look ahead somewhat to a possible future time at which a purely automatic combination of digital and analog methods can be effected. Many problems will have to be solved to do this - such as design of high-speed digital-to-analog and analog-to-digital translators - before such a combination will become possible. In such a machine, the analog components would be used wherever their other advantages outweighed the need for precision.

This prospect is, however, something very definitely in the future, and until this means or others become practically realizable, the advantages of the analog type of computer for control system work will, in almost all cases overbalance the higher accuracy of the digital computer.

CHAPTER VIII

THE ANALOG COMPUTER

SECTION 1 - INTRODUCTION

In the preceding chapter, it was shown that the operational amplifier analog computer is the most valuable device now in existence for the machine solution of control systems dynamical problems. The present chapter considers material relevant to the application of this type of automatic computation to such problems.

The first section deals with d.c. amplifiers and explains how they are used to perform certain mathematical operations required for the solution of differential equations. In this section it is assumed that the amplifier is an ideal device. Hence, this discussion will be completely applicable only in those cases where the various approximations made in idealizing the operational amplifier introduce negligible error.

The second section will consider real amplifiers (i.e. non-ideal) and show how and why the results obtained from such a real amplifier differ from the results obtained from an ideal amplifier. In particular, the effects of non-infinite gain in the amplifier, of drift, and of grid current will be discussed.

The third section will consider in more detail the use of the operational amplifier as a summer, a sign changer, and as an integrator.

This will be followed by material related to the interconnection and use of operational amplifiers in the solution of single and simultaneous linear differential equations with constant coefficients.

In this part of the chapter, the manner in which a com-

puter is set up for a representative set of suitable equations will be considered. This will include the determination of "gains," the selection of appropriate voltage levels, the computation of suitable resistor values and potentiometer settings, and the arrangement of the apparatus for maximum utilization of equipment and minimization of the possibility of misbehavior.

The concluding section of this chapter will then take up the important practical question of the representation in machine form of non-linearities. As pointed out in chapter II any real control system will contain non-linear effects such as backlash, Coulomb friction, and threshold and very often, excessive efforts may be made to minimize their effect because of ignorance of their influence on system behavior. Other non-linearities, such as spring pre-loading or varying spring constants, may actually be built into a system in order to produce a dynamic effect which cannot be obtained by linear means.

The means of representing threshold, acceleration, velocity, or displacement limitations, backlash, Coulomb friction, spring preload, and varying spring constants will be discussed. Certain limitations on the means of representing these effects will be pointed out.

The non-linear problem is one in which computer operation has a very great advantage over other methods for determining the response, since analytic solution of equations with non-linearities is inherently a tedious step-by-step process. The automatic computational methods seem to offer the only hopeful technique of dealing with these very important properties of the system.

SECTION 2 - "IDEAL" D.C. AMPLIFIER OPERATION

This section discusses ways in which operational amplifiers may be used to perform certain mathematical processes used in the solution of differential equations. These processes usually consist of algebraic addition, multiplication, integration, and differentiation. In the operational amplifier, these quantities are voltages which are analogous to the dependent variables of the differential equations written for the systems to be studied.

In general, the operational amplifier consists of a high gain dc voltage amplifier, an input impedance, and a feedback impedance. Figure VIII-1 illustrates the manner in which these elements are interconnected.

The impedance, Z_i , is added to this figure to represent the input grid resistor of the voltage amplifier.

The overall gain of this amplifier may be found by writing the four equations (VIII-1) through (VIII-4), and

$$(VIII-1) \quad e_i = i_i Z_i + e_g$$

$$(VIII-2) \quad e_o = i_i Z_f + e_g$$

$$(VIII-3) \quad e_g = -K e_o$$

- * Operational amplifier design requirements necessitate an odd number of amplifier stages. Hence, the minus sign.

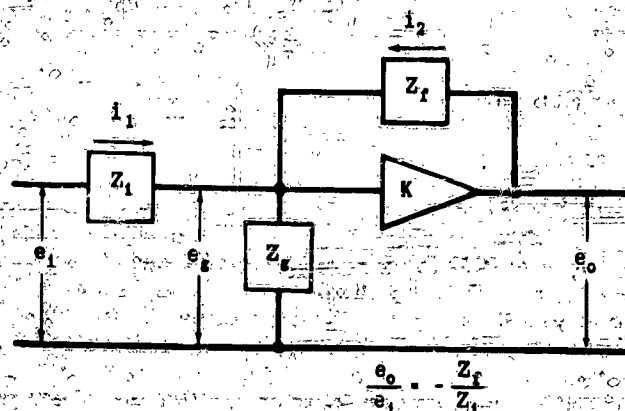


Figure VIII-1. Operational Amplifier.

$$(VIII-4) \quad e_2 = (i_1 + i_2) Z_f$$

solving for the ratio e_0/e_1 . This overall gain, as found from these equations, is given by

$$(VIII-5) \quad \frac{e_0}{e_1} = \frac{Z_f}{Z_i} \frac{1}{1 + \frac{1}{K} \left(1 + \frac{Z_f}{Z_i} \right)}$$

In this section, the amplifier will be considered an ideal device. As shown by equation (VIII-5), a very high gain dc voltage amplifier will make the operational amplifier relatively independent of all quantities except the ratio of feedback impedance (Z_f) to input impedance (Z_i). An ideal operational amplifier then, would have a dc voltage amplifier of infinite gain. For this ideal case, the relationship between the input voltage and the output voltage is given by equation (VIII-6), where the minus sign indicates a change in polarity.

$$(VIII-6) \quad e_0 = -\frac{Z_f}{Z_i} e_1$$

The ratio $-Z_f/Z_i$ may be thought of as a multiplying factor or operator by which the input voltage is multiplied to get the output voltage. In particular, if these two impedances are resistors, the input voltage will be multiplied by the ratio of these resistors, and will have its sign changed. The operation performed by the operational amplifier is that of multiplication by a constant and sign changing for this case. If the ratio is unity, the operational amplifier serves only as a sign changer.

If the feedback impedance is a capacitive reactance, i.e.: $Z_f = 1/(sC)$; and the input impedance, a resistance, the operator becomes

$$(VIII-7) \quad -\frac{Z_f}{Z_i} = -\frac{1}{R_i C s}$$

The relation between the input voltage and output voltage is given by equation (VIII-8).

$$(VIII-8) \quad e_0 = -\frac{1}{R_i C s} e_1$$

The equivalent relation in the time domain is

$$(VIII-9) \quad e_0 = -\frac{1}{RC} \int e_1 dt$$

For this case, the operation performed by the operational amplifier is that of integration, sign changing, and multiplication by a constant ($1/R_i C$).

The operational amplifier may be made to differentiate by interchanging the resistance and capacitive reactance used above so that the relation between the input and output voltages will be as given by equation (VIII-10).

$$(VIII-10) \quad e_0 = -R_f C s e_1$$

The equivalent relation in the time domain is

$$(VIII-11) \quad e_0 = -R_f C \frac{d(e_1)}{dt}$$

In this case, the operational amplifier will multiply this derivative by $R_f C$ and change its sign. As will be pointed out in section VIII-3, the use of the operational amplifier as a differentiator may not be desirable because of the presence of noise.

If the feedback impedance of the operational amplifier consists of a parallel combination of resistance and capacitive reactance so that $Z_f = R_f / (R_f C s + 1)$, the resulting relationship between the input voltage and the output voltage is given by equation (VIII-12).

$$(VIII-12) \quad e_0 = -\frac{R_f}{R_i} \left(\frac{1}{R_f C s + 1} \right) e_1$$

In the time domain, this relation is given by equation (VIII-13).

$$(VIII-13) \quad e_0(t) = \mathcal{L}^{-1} \left[\frac{R_f}{R_i} \left(\frac{e_1(s)}{R_f C s + 1} \right) \right]$$

For a unit step input, the output is a first-order lag as shown in (VIII-14).

$$(VIII-14) \quad e_0(t) = -\frac{1}{R_i C_f} \left(1 - e^{-\frac{t}{R_f C_f}} \right)$$

Before discussing the process of addition with the operational amplifier, it is of importance to note that the grid voltage, e_2 , of the dc voltage amplifier remains relatively unchanged while the input is varied. That is, as this grid voltage tends to change, the output voltage tends to oppose this change. To show that this is true, equations (VIII-1) through (VIII-4) are solved for the ratio of grid voltage to input voltage:

$$(VIII-15) \quad \frac{e_2}{e_1} = \frac{1}{K \frac{Z_i}{Z_f} + \left(1 + \frac{Z_i}{Z_f} \right)}$$

It is apparent that, for gains approaching infinity, as does that of the ideal amplifier, the grid voltage remains at zero.

Figure VIII-2 illustrates the manner in which the elements are interconnected to perform the process of addition. The sum of the currents flowing into point A must equal the sum of the currents flowing out of point A:

$$(VIII-16) \quad \frac{e_1 - e_A}{Z_1} + \frac{e_2 - e_A}{Z_2} + \frac{e_3 - e_A}{Z_3} + \frac{e_0 - e_A}{Z_f} = \frac{e_A}{Z_s}$$

The grid current is assumed zero in the voltage amplifier and does not enter into equation (VIII-16). As shown in the preceding paragraph, the grid voltage is zero in the ideal operational amplifier. Equation (VIII-

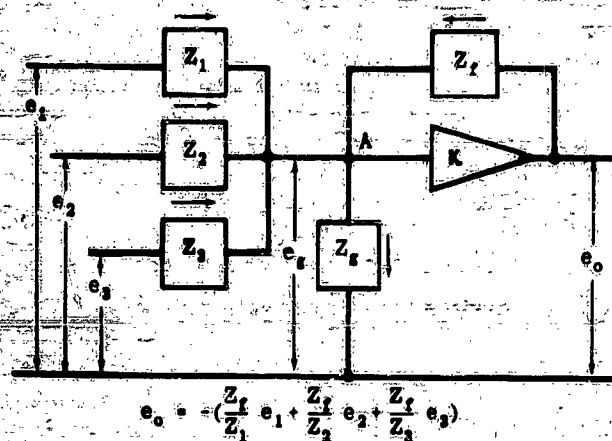


Figure VIII-2. Addition

SECTION 3 - NON-IDEAL OPERATIONAL AMPLIFIERS

In Section VIII-2 it was shown how the operational amplifier, if ideal, could be used to perform certain mathematical operations. Actually, the operational amplifier is not ideal. The dc voltage amplifier does not have infinite gain, and is subject to "drift," grid current, electrical noise, and limiting. In addition to these undesirable characteristics in the dc voltage amplifier, the other components which make up the operational amplifier may not have the properties assumed. In this section the effect of these non-ideal properties on the operations performed will be discussed qualitatively. Quantitative material will be developed in section VIII-4.

The dc voltage amplifiers used in electronic analog computers usually have voltage gains ranging in the order of 5,000 to 100,000. As is apparent from equation (VIII-5), the fact that the gain is not infinite puts certain restrictions upon the permissible values of input and feedback impedances. These restrictions will be further discussed in section VIII-4.

When the input to a high gain and properly balanced dc voltage amplifier is zero the output should also be at zero potential. However, there may be small current changes (drifts) in the tubes. These drifts may be caused by slow changes in dc supply potentials, in cathode emission due to filament supply changes, and in the resistance of resistors due to changes in the ambient temperature. When these current changes occur in the first tube of the dc voltage amplifier, their effect will be multiplied by the gain of the amplifier with a resulting large change in the output. When such an amplifier is used as a computer, the output signal will then have this drift voltage added to that calculated for the ideal operational amplifier.

The dc voltage amplifiers in operational amplifiers are used in what is termed "class A operation." That is, the grid is always negative with respect to the cathode.

16) then takes the form of (VIII-17)

$$(VIII-17) \quad \frac{e_1}{Z_1} + \frac{e_2}{Z_2} + \frac{e_3}{Z_3} = -\frac{e_o}{Z_f}$$

which, when rearranged can be written in the form

$$(VIII-18) \quad e_o = -\left[\frac{Z_f}{Z_1} e_1 + \frac{Z_f}{Z_2} e_2 + \frac{Z_f}{Z_3} e_3\right]$$

From this latter equation it can be seen that the output of the operational amplifier is the sum of the inputs multiplied by their respective operators. Though only three inputs were considered above, any number can be used. The operation of integration or differentiation, multiplication by constants, and summing can be combined in a single operational amplifier.

In this section, it has been shown how the operational amplifier may be used to perform certain mathematical processes needed for the solution of linear differential equations with constant coefficients. The manner in which non-linearities may be included with the operational amplifier will be discussed in section six of this chapter.

When the grid is negative with respect to the cathode, a current flows away from the grid terminal, as shown in figure VIII-3. This current is of the order of micro-amperes; but, when flowing through a sufficiently high impedance, it may cause an appreciable voltage drop across that impedance. As was shown in the previous section, the input grid voltage of the operational amplifier is effectively at zero potential for any reasonably large gain in the dc voltage amplifier. For this reason, the grid impedance is not shown in figure VIII-3.

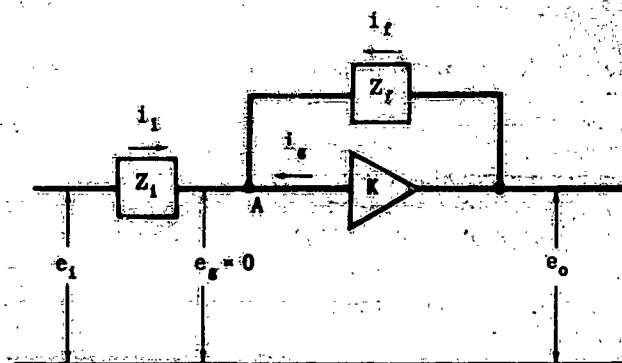


Figure VIII-3. Operational Amplifier with Grid Current

In figure VIII-3, the sum of the currents entering point A must be zero as shown by equation (VIII-19).

$$(VIII-19) \quad \frac{e_1}{Z_1} + \frac{e_2}{Z_2} + i_g = 0$$

This equation, when solved for the output voltage, becomes

$$(VIII-20) \quad e_o = -\left(\frac{Z_f}{Z_1} e_1 + i_g Z_f\right)$$

Equation (VIII-20) shows that the output voltage consists of the component which would be obtained from an ideal

operational amplifier, plus an additional component determined by the magnitude of the grid current and the feedback impedance.

Generally speaking, noise may be defined as any undesirable voltage which may appear with the desired voltage. The average value of the noise has been discussed in connection with drift and grid current. In the following discussion, noise will refer to that component which has a zero average value. The sources of noise are many. Within the dc voltage amplifier, the two primary sources of noise are that generated in resistors when currents are flowing through them, and the thermal noise generated within the vacuum tubes due to the uneven "boiling off" of the electrons from the cathodes. Of greater importance is, perhaps, the noise due to "pickup." The wires used to interconnect the operational amplifiers cannot be completely shielded, and will pick up stray fields radiated by other components used in the analog computer. The presence of noise in these computers limits the lowest permissible voltage which can be operated upon by the operational amplifiers without having the noise appear as an appreciable part of the final result.

It was mentioned in the previous section that differentiation with the operational amplifier may not be recommended due to the presence of noise. If the noise voltage has small rise and decay times, as do the thermal and resistor types, the operational amplifier connected as a differentiator will differentiate these voltages. If the input voltage is varying slowly, it is easy to see that the output signal may become largely noise. In addition the differentiated noise amplitude may become so great that amplifiers overload.

The dc voltage amplifier used in the operational amplifier is limited to a maximum voltage which can be obtained at the output. Above this value, the output stage saturates. To help eliminate the possibility of unknowingly exceeding the linear operation range of the amplifiers, most electronic analog computers make use of indicators which signal when limiting occurs.

As was pointed out in the previous section, the operation performed by the operational amplifier depends primarily upon the input and feedback impedances. The operational amplifier may be considered as non-ideal in the sense that these impedances may not be as assumed. While the values of resistors and condensers used for these impedances are stated, the fact that they may differ by whatever tolerance has been assigned must be kept in mind. In addition to the values, it is of importance to know if the components are effectively pure. That is, if a resistor contains only resistance and a capacitor only capacitive reactance. At the frequencies used in analog computers the inductive effects of resistors are negligible and need not be considered. Capacitors, however, have a certain amount of leakage through them so that an accurate representation of a capacitor would show a resistor in parallel to represent the leakage. To reduce this leakage to a negligible value, highest quality condensers

are used in the better analog computers.

As will be shown in section five, the gains used in the operational amplifier may correspond to the coefficients of the differential equation being analyzed. Since there is an infinite number of possible coefficients, an infinite number of possible gains could be required. It is not practical to obtain these gain variations by varying the feedback impedance, since when integrating, this would mean a very large variable capacitor would have to be employed. For the same reason it is not very practical to vary the input impedance appreciably. Furthermore, if small values of input impedance were required, the sources of input voltage would be excessively loaded.

To make possible the desired gain variations without varying the input and feedback impedances, electronic analog computers have provisions for inserting a potentiometer in the input of the operational amplifier as shown in figure VIII-4.

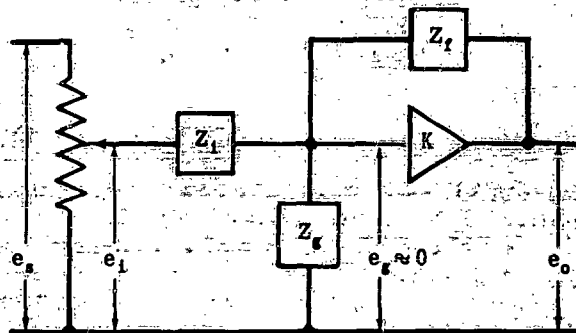


Figure VIII-4. Operational Amplifier with Potentiometer

In effect the potentiometer may be considered as an amplifier with a gain variable from zero to unity. As may be noted in the figure, the fact that the grid voltage is always near zero effectively places the input impedance across the potentiometer as illustrated in figure VIII-5.

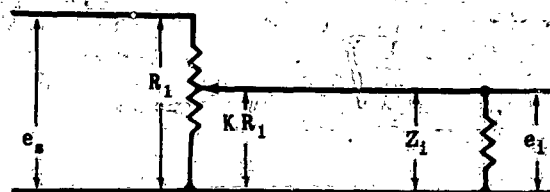


Figure VIII-5. Equivalent Circuit

If KR_1 is that fraction of the potentiometer R_p across Z_i , the overall gain (e_1/e_s) is

$$(VIII-21) \quad \frac{e_1}{e_s} = K \left[\frac{1}{K \frac{R_1}{Z_i} (1-K) + 1} \right]$$

From this expression, it is apparent that the gain ahead of the operational amplifier is the fraction, K , of potentiometer setting only when the input impedance is much larger than the potentiometer resistance. Equation (VIII-21) may be used to construct a family of curves relating the fractional potentiometer setting to the actual potentiometer gain for various values of input resistances.

SECTION 4 - EFFECTS OF NON-IDEAL OPERATION

In section three, it was pointed out that operational amplifiers are not ideal devices. They may be assumed ideal, however, if the effects which make the amplifiers differ from ideal can be made negligible. In this section, measures which minimize these effects will be discussed so that the operational amplifier may be used as ideal with a negligible error in results. If such steps were not taken, it would be necessary to make time consuming calibration and correction computations.

In the dc voltage amplifiers used with most analog computers, the grid impedance, Z_i , is made much larger than the input and feedback impedances. For these cases, equation (VIII-5) may be simplified to that of (VIII-23).

$$(VIII-22) \quad \frac{e_o}{e_i} = \frac{Z_f}{Z_i} \frac{1}{1 + \frac{1}{K} \left(1 + \frac{Z_f}{Z_i} \right)}$$

For purposes of discussion, equation (VIII-23) may be put in the form of (VIII-23)

$$(VIII-23) \quad \frac{e_o}{e_i} = \frac{Z_f}{Z_i} \left| \frac{K}{K+1} \right| \left| \frac{1}{1 + \frac{Z_f}{Z_i} \frac{1}{K+1}} \right|$$

so that the error terms of the ratio $-Z_f/Z_i$ will consist of a term which is a function only of the gain, and a second term which is a function of both gain and the impedances.

For any given operational amplifier, the first bracketed term of equation (VIII-23) is a constant, and usually cannot be changed. If the error introduced by this term is negligible in itself, attention can be given the second term. The second term will be frequency sensitive if the operation performed is integration or differentiation, and will introduce a phase shift. It is assumed that the permissible error is known, and that a problem exists of determining an operating procedure which will keep the gain and phase changes within the acceptable limits. An examination of the second term will determine these limits of operation.

Consider the problem of addition. It can be shown that the equation corresponding to (VIII-23) is (VIII-24),

$$(VIII-24) \quad e_o = \left(\frac{R_f}{R_1} e_1 + \frac{R_f}{R_2} e_2 + \frac{R_f}{R_3} e_3 \right) \left(\frac{K}{K+1} \right) \left[\frac{1}{1 + \frac{R_f}{R_t} \frac{1}{(K+1)}} \right]$$

where R_t is the parallel combination of the input impedances. From the third, bracketed term, it is apparent that there is a limitation as to the number of voltages which may be added for any given feedback resistor, gain, and acceptable error in the non-ideal operational amplifier.

In the case of differentiation, the second bracketed term of (VIII-23) becomes that of (VIII-25).

$$(VIII-25) \quad \frac{1}{1 + sC_f R_t \frac{1}{K+1}}$$

If $j\omega$ is substituted for s , (VIII-25) may be separated into the amplitude and phase expressions of (VIII-26) and (VIII-27).

$$(VIII-26) \quad \text{Amplitude} = \frac{1}{\sqrt{1 + \left(\frac{\omega C_f R_t}{K+1} \right)^2}}$$

$$(VIII-27) \quad \phi = -\tan^{-1} \left(\frac{\omega C_f R_t}{K+1} \right)$$

Equation (VIII-26) when combined with the first bracketed term of (VIII-23) provides the means for determining the permissible range in component values and frequency which will keep the amplitude error within the acceptable limits. Equation (VIII-27) provides the same information necessary for keeping the phase shift within acceptable limits.

When integration is performed by the operational amplifier, the amplitude and phase contributed by the second bracketed term of equation (VIII-23) are those of (VIII-28) and (VIII-29).

$$(VIII-28) \quad \text{Amplitude} = \frac{\omega R_f C_f (K+1)}{\sqrt{1 + \omega R_f C_f (K+1)^2}}$$

$$(VIII-29) \quad \phi = \tan^{-1} \frac{1}{\omega R_f C_f (K+1)}$$

As with the case of differentiation, these equations provide the means for determining the permissible range in component values and frequency which will keep the amplitude and phase errors below acceptable values.

To discuss the steady-state accuracy of operational amplifiers when used as integrators, it is convenient to do so in terms of the response to a unit step input. The unit step is chosen because the error will be larger than for any other possible input.

Where the input impedance is a resistor and the feedback impedance a capacitor, equation (VIII-23) may be written in the form of (VIII-30).

$$(VIII-30) \quad \frac{e_o}{e_i} = - \frac{K}{sC_f R_t (K+1) + 1}$$

It may be noticed that this equation is of the same form as (VIII-12). That is, the non-ideal operational amplifier behaves like the ideal amplifier with a very large time lag.

The response of the integrating operational amplifier to a unit step input, as found from (VIII-30), is described by (VIII-31).

$$(VIII-31) \quad e_o(t) = -K \left(1 - e^{-\frac{t}{(K+1)R_t C_f}} \right) \text{ for } (K+1)R_t C_f \sim KR_t C_f$$

When expanded in a power series, (VIII-31) takes the form of

$$(VIII-32) \quad e_o(t) = -K \left(\frac{t}{(K+1)R_t C_f} - \frac{t^2}{2(K+1)^2 R_t^2 C_f^2} + \frac{t^3}{6(K+1)^3 R_t^3 C_f^3} - \dots \right)$$

or

$$(VIII-33) \quad e_o(t) = -\frac{t}{R_1 C_f} \left(1 - \frac{t}{2\tau}\right)$$

for $t/\tau \ll 1$. Since the ideal integrator would have an output of $t/R_1 C_f$, the term $-t/2\tau$ is that fraction by which the output of a non-ideal integrator is in error for $t < \tau$. As apparent from (VIII-33), this error may be made small by making τ large. Since the $R_1 C_f$ product is a constant for a given operation, the error may be made small by using operational amplifiers with a large value of K .

It can be shown that the output voltage of an operational amplifier with drift is

$$(VIII-34) \quad e_o = \frac{Z_f}{Z_i} \frac{1}{1 + \frac{1}{K} \left(1 + \frac{Z_f}{Z_i}\right)} e_i + \frac{1}{K + 1 + \frac{Z_f}{Z_i}} \left[\frac{Z_f}{Z_i} + 1\right] e_d$$

In this equation, the term e_d is the drift voltage measured at the output of the dc voltage amplifier when the input is grounded. From (VIII-34), it is seen that the output of the operational amplifier consists of a term which would exist in a driftless amplifier, and a term which exists due to the drift voltage. For large values of gain, this second term reduces to $(1/K)(Z_f/Z_i) + 1 e_d$. While it may appear that a sufficiently large gain would eliminate the effect of drift, it must be remembered that the drift which appears in the output is due to the drift in the first stage of the dc amplifier. Any increase in gain will also increase the

drift voltage.

If the operation to be performed is that of integration, the drift term takes the form $(1/K) [1/(sR_1 C_f) + 1] e_d$. That is, the drift voltage is integrated as well as the input voltage.

Since the error due to drift is a function of time, it is a simple matter to measure the drift of an operational amplifier with no input and note the time it takes to exceed a specified negligible value. Any solution extending beyond this time is not within the required accuracy. The error due to drift may be minimized by using large amplitude signals, since the drift voltage is then only a small proportion of the total output signal.

As shown by (VIII-20), the effect of grid current is to add a voltage to the output which is a product of the grid current and the feedback impedance. Though this grid current is small in well-designed amplifiers, a sufficiently large value of feedback impedance may add an appreciable voltage to the output. The effect of grid current may usually be made negligible by restricting the feedback impedance to values of one megohm or less.

In the next section, the use of operational amplifiers for the solution of differential equations will be discussed. Though the assumption that the amplifiers are ideal will be made, it will be noticed that the value of the components and voltage levels used tend to minimize the undesirable effects of non-ideal operational amplifiers.

SECTION 5 - SOLUTION OF DIFFERENTIAL EQUATIONS

This section will discuss in detail the process of interconnecting operational amplifiers for the solution of a set of simultaneous linear differential equations with constant coefficients. The essentials of the process involved shall be illustrated by use of an example involving a fairly complex set of equations. Specifically, the equations will be those of the longitudinal motion of an aircraft under the control of a simple autopilot. The autopilot is assumed to consist of a second-order servo motor, a perfect rate plus displacement equalizer, and a single time lag hydraulic system between the servo motor and the control surface deflection.

The equations of motion in question are as follows

$$(VIII-35) \quad \ddot{u} = -0.009\ddot{u} - 0.0362\dot{u} - 32.2\dot{u}$$

$$(VIII-36) \quad \ddot{w} = -0.0621\ddot{w} - 0.612\dot{w} + 807\theta + 25.2\theta$$

Airframe Equations

$$(VIII-37) \quad \ddot{\theta} = -0.0006\ddot{\theta} - 0.0028\dot{\theta} - 1.03\theta + 8.97\delta$$

$$(VIII-38) \quad \ddot{\sigma} = -12\ddot{\sigma} - 400\dot{\sigma} + K_1(\dot{\theta} + K_2\theta)$$

Servomotor Equation

$$(VIII-39) \quad \delta = \frac{K_3 \sigma}{0.1s + 1} + \delta_1$$

Control Surface Equation

where: (All quantities are "perturbed" values from a steady-state equilibrium condition.)

- u is the forward velocity of the aircraft.
- w is the vertical velocity of the aircraft.
- θ is the pitch angle of the aircraft.
- σ is the output rotation of the servo motor.
- δ is the total control surface deflection.
- δ_1 is an input disturbance to the surface.
- K_1 & K_2 are equalizer gains.
- K_3 is the gain between servo motor and surface.

It will be noticed that the three constants K_1 , K_2 , K_3 , have not been specified. This has been done to provide a simple illustration for later use, to show how the analog computer may be used for synthesis. With these parameters available for change it is possible to optimize the performance of the physical system for which the balance of parameters is fixed.

It should be noticed that the equations are written so that the highest derivative of the principal variable in each equation is isolated on the left side. This is done to facilitate setting up the computer.

The control surface equation is an exception to this rule. However, the equation contains a simple time lag, and, as was stated in section VIII-2, this is most conveniently represented by an amplifier with a parallel RC network in the feedback.

Consider the meaning of the first of the airframe equations. Expressed in words, it says that the rate of change of the quantity u is equal to a linear combination of u , w , and θ . In terms of the apparatus, if u , w , and θ are each multiplied by the appropriate constant and added together (algebraically), the result is \dot{u} .

Section VIII-2 has already demonstrated that these operations can be carried out by operational amplifiers. However, close scrutiny shows that the quantity \dot{u} is needed nowhere in the solution. What is really wanted is rather the quantity u . As has been mentioned before in this chapter, summing and integration may be combined in a single amplifier. The first equation may then be set up by the following arrangement of apparatus:

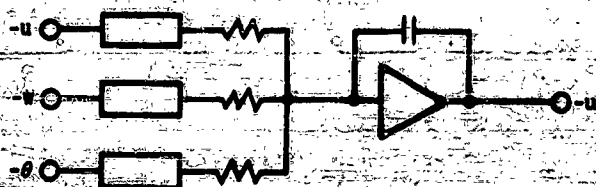


Figure VIII-6. Mechanization of the \dot{u} equation

It should be noted that this connection of apparatus is equivalent to writing (VIII-35) in the form:

$$(VIII-40) \quad \int \dot{u} dt = u = \int (-0.009u - 0.0362w - 32.2\theta) dt$$

Nothing has yet been said about where w and θ come from. However, it should be clear that by mechanizing the other equations and making the proper connections so that each quantity appearing as the output of an amplifier is fed to the appropriate input, the various quantities involved will be related as described by the equations. Hence, if the completely mechanized system is disturbed in any manner at all, and thus caused to "move," its motion can only be that which satisfies the equations.

Another thing to be observed about figure VIII-6 is that while potentiometers, input resistors, and a feedback capacitor are indicated by the diagram, their numerical values do not appear. This will be taken care of later. They will be chosen so that not merely the general functional form of the equation is realized, but also so that the coefficients have the magnitudes occurring in the equations.

The second of the airframe equations, (VIII-36), may similarly be symbolically represented in terms of the apparatus as:



Figure VIII-7. Mechanization of the \dot{w} equation

Here it was not possible to carry out the summation and the integration in a single amplifier, since \dot{u} is needed for use in the equation for $\dot{\theta}$. Therefore the input quantities had to be summed separately. Also, the two reversals of sign in the summer and the integrator gives $+w$ instead of $-w$ which is needed in the $\dot{\theta}$ equation. Therefore another amplifier has to be used to obtain the correct algebraic sign for this quantity.

The third aircraft equation of motion may be set up in much the same way. The principal difference is that two integrations are required.



Figure VIII-8. Mechanization of the $\dot{\theta}$ equation

Here it will be observed that it has been necessary to provide both algebraic signs of θ and of $\dot{\theta}$. The negative signs are used for the $\dot{\theta}$ and \dot{u} equations but the positive ones will be required for the relation defining $\dot{\theta}$. The $+\dot{\theta}$ amplifier was taken off in a side branch instead of cascading all four amplifiers to reduce the number of amplifiers in the chain.

The $\dot{\sigma}$ and $\dot{\delta}$ equations may now be set up, and the interconnections of the apparatus are as shown by figure VIII-9.

Notice that $+\sigma$ is not required since the additional sign change in the $\dot{\delta}$ amplifier gives rise to $+\delta$, which is needed for feeding back into the airframe equations.

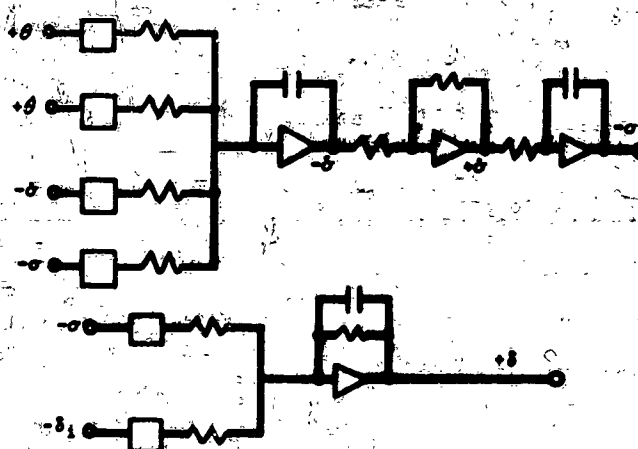


Figure VIII-9. Mechanization of $\dot{\sigma}$ and $\dot{\delta}$ Equations

Figure VIII-9 represents the functional form of the above equations. There still remains the task of putting the appropriate numbers into the mechanization.

First it is necessary to select appropriate voltage levels, (i. e., how many volts are to correspond to

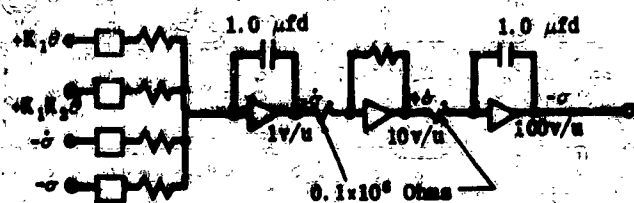


Figure VIII-11. Mechanisation of Servo Motor Equation

By hypothesis, K_1 and K_2 are not known, since this is assumed to be a synthesis problem which requires finding values for suitable response of the system. This means that these two coefficients would have to be determined experimentally by trial of various combinations of resistors and potentiometer settings and observation of the transients occurring in response to various inputs of interest. In this, the actual gains K_1 , K_2 and K_3 would then have to be calculated from the

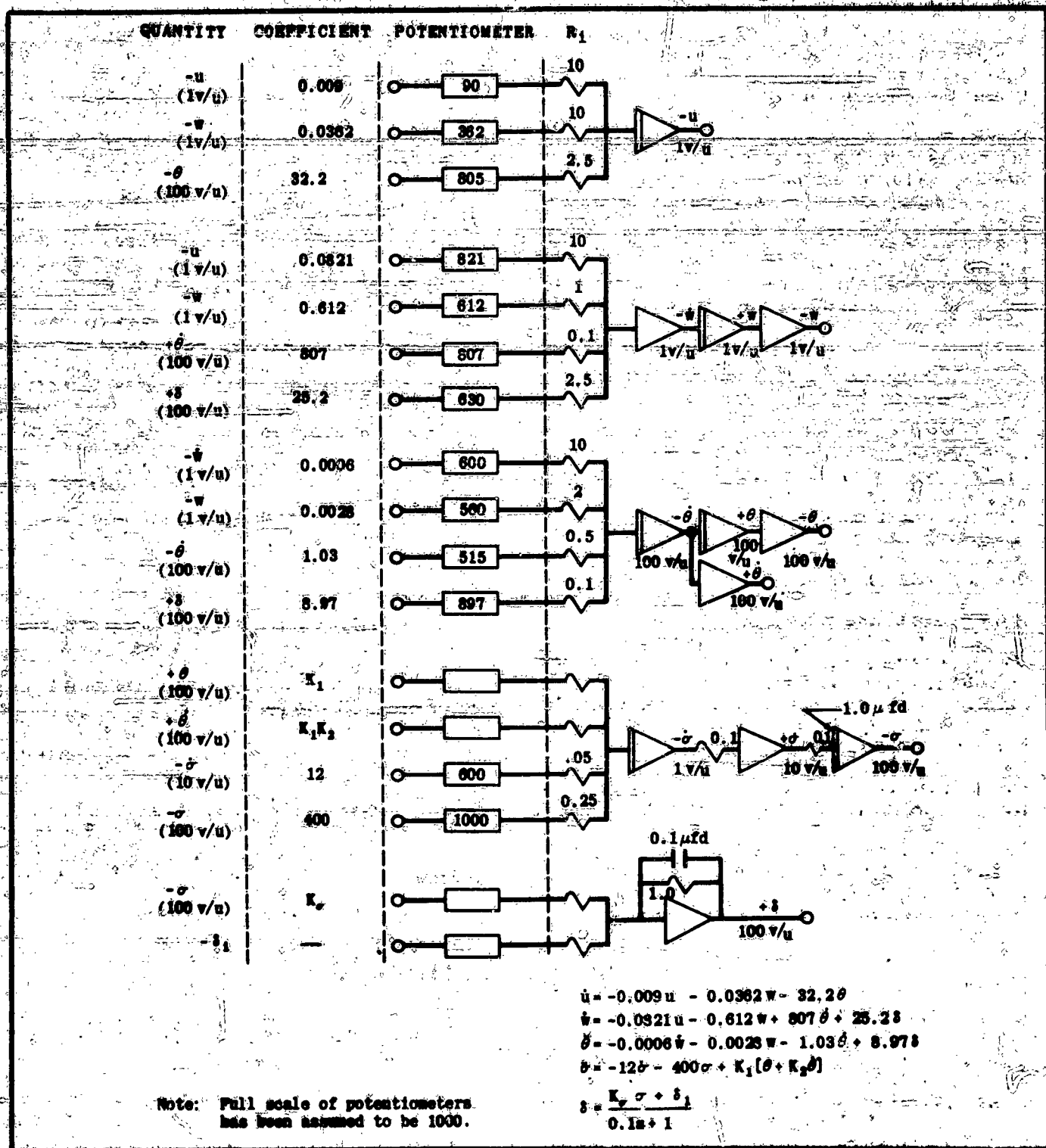


Figure VIII-12. Analog Computer Setup

Chapter VIII Section 6

R_1 's and the potentiometer settings, after these were known.

However, there is another problem. The gain from σ to δ is 400; that is, $\partial\delta/\partial\sigma = 400$; and while gains of 100 or so across a single operational amplifier are practicable if the amplifier is a good one, 400 is definitely much too high. This disposes of the obvious solution of making $C = 10^{-8}$ farads, $R_{12} = 2500$ ohms so that $1/R_{12}C = 400$.

This gain does not all have to be taken at the input σ to the circuitry, however. It makes no difference if part of the gain is taken there, and the rest in passing through the other amplifiers in the cascaded chain.

Suppose $C = 10^{-6}$, $R_{12} = 1/4 \times 10^6 = 250,000$ ohms; then a gain of 4 is picked up across the first integrator. If then the input resistor to the sign changer is made 0.1×10^6 ohms, and its feedback resistor a megohm, a voltage gain of ten results at this point. A further gain of 10 may be taken in the last integrator by using a 0.1 megohm input resistor and a microfarad feedback capacitance. (This is preferable to a megohm and a 0.1 microfarad, since the first method reduces the impedance level.)

Then the gain around the closed σ loop (without velocity feedback) is 400 as it should be.

It is of great importance to notice, however, that something essentially different has been done here than if the gain were all taken at the input to the first integrator. In the latter case, the voltage level from input to output of the σ amplifier would not have been changed. In the way the gain was actually taken, the voltage level has been raised twice, as shown in the diagram of figure VIII-11.

If one volt is placed at the input of the sign changer, for example, and this represents one unit of σ , then the ten volts resulting at the output of that amplifier can still represent only that one unit of the physical quantity δ ; and hence the voltage level, the number of volts/unit of σ , has been changed. It is essential to bear this concept in mind when the voltage representing σ is used as in the δ equation here, as it is necessary to know the σ voltage level so that the gain into the next equation may be properly calculated.

The gain $\partial\delta/\partial\sigma$ requires no particular effort; a gain of 20 may be taken at this input by using a 50,000 ohm resistor, and setting the associated potentiometer for 600 (uncorrected).

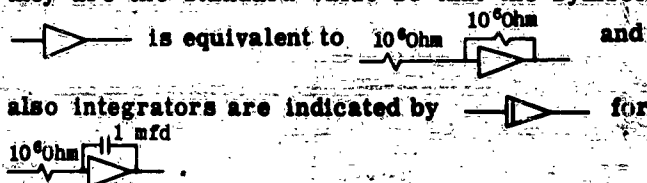
The gain $\partial\delta/\partial\sigma = K_1$ in the next equation is again not known. The voltage level of δ may arbitrarily be taken as 1,000 volt per unit which is convenient for feeding δ into the \dot{w} and θ equations, and K_1 ultimately calculated on this basis.

It remains only to determine the feedback resistor and capacitor in the δ equation so that the time lag has the right magnitude to complete the mechanization of the equations (VIII-35) through (VIII-39).

As was shown in section VIII-2 the RC product here must be the value of the lag, 0.1 sec. in this case. This is easily obtained by using, say, a megohm resistor, and a 0.1 mfd capacitor. As has been previously remarked, this can be a good paper condenser; the high-quality polystyrene type needed for integration is not required for use in time lag circuitry.

Figure VIII-12 shows the complete diagram for the mechanization of the equations under consideration. All numbers are included except those to be determined by the synthesis; those not computed above may easily be verified if it is desired to do so.

This diagram has been drawn on the basis that all resistances are one megohm and all capacitances one mfd unless it is otherwise explicitly indicated. The convention has also been adopted that input and feedback resistances of a sign changer are not drawn if they are the standard value so that the symbol



As a synthesis problem, the one here considered is very simple. However, if it were desired to do more than merely fix certain gains, that is, if it were desired to investigate the effects of various additional lags or to introduce other types of control, it is easy to see that these can be introduced into the computer by introducing new apparatus and interconnecting it properly.

Another interesting and important aspect of such problems is to determine the effects of certain non-linearities, such as Coulomb friction or limiting, upon the behavior of the system. The machine means of representing such discontinuities will be the subject of the next and final section of this chapter.

SECTION 6 - SIMULATION OF NON-LINEARITIES

(a) INTRODUCTION

In the previous sections of this chapter, it was shown how the electronic operational amplifier could be used to obtain the solution of linear differential equations with constant coefficients. As was pointed out in chapter VI, real physical systems may not be accurately described by this type of equation because of non-linearities which may exist. This section will discuss

methods by which it is possible to introduce these non-linearities in the analog computer for more accurate representation of physical systems.

Since the simulation of non-linearities in the analog computer makes frequent use of the diode or polarized relay, a brief discussion of the general properties of these devices will aid the understanding of the circuits to be discussed.

The ideal diode is a device which will behave like a resistor of zero resistance for currents flowing in one direction, and a resistor of infinite resistance for currents flowing in the opposite direction. The non-ideal diodes possess neither of these features. In figure VIII-13 typical characteristic curves are shown for non-ideal diodes. The reciprocal of the slopes of these curves have the units of resistance. As shown by the relatively small slope for negative values of e , the resistance of the diode (R_d) is quite large in value. As e increases through zero and assumes positive values, the slope of the curve increases, and the corresponding resistance decreases. For still larger values of e , the slope becomes relatively constant and determines the lowest resistance a diode may have.

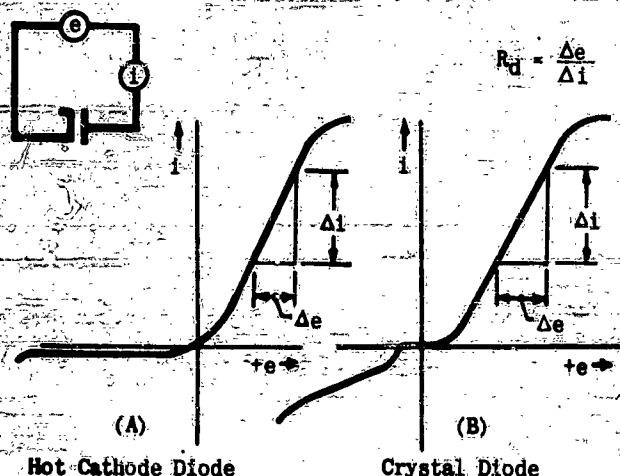


Figure VIII-13. Diode Characteristics

The polarized relays, as used in computation, are very sensitive devices. When no voltage is applied to the electromagnet, the armature is in a neutral position. When a voltage is applied to the electromagnet, the armature closes either of two possible contacts depending upon the polarity of the voltage applied.

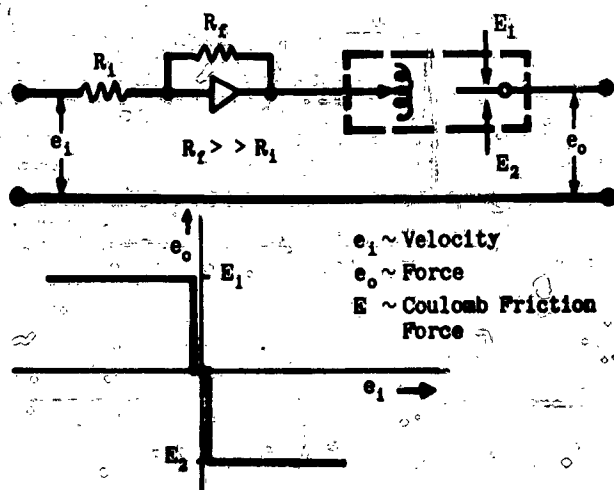


Figure VIII-14. Simulation of Coulomb Friction with Polarized Relay

(b) COULOMB FRICTION

To simulate coulomb friction, it is necessary to have a device which will accept a voltage (representing velocity) and deliver a constant voltage (representing force), which will reverse in polarity when the voltage representing velocity is reversed.

The polarized relay could be used alone to do this, but the voltage representing velocity would have to be relatively large before the contacts would close in one direction or the other. To overcome this difficulty, a high gain amplifier is used in conjunction with the relay as shown in figure VIII-14. With the aid of this amplifier, the relay contacts close at a very small value of input voltage. The voltages E_1 and E_2 applied to the relay contacts are the voltages which represent the coulomb friction force.

In figure VIII-15 the representation of coulomb friction is accomplished through use of biased diodes and operational amplifiers.

For zero output voltage the diodes are biased by the battery (B) so that both diodes represent very high resistances and effectively open the circuit. Since the battery (B) and the resistors (R) form a dc bridge, the voltage at point a equals the voltage at point b. For an output voltage equal to one-half the battery voltage, it may be noticed that one of the diodes must have twice the bias it originally had, while the other has zero bias. Any increase of e_o above this value

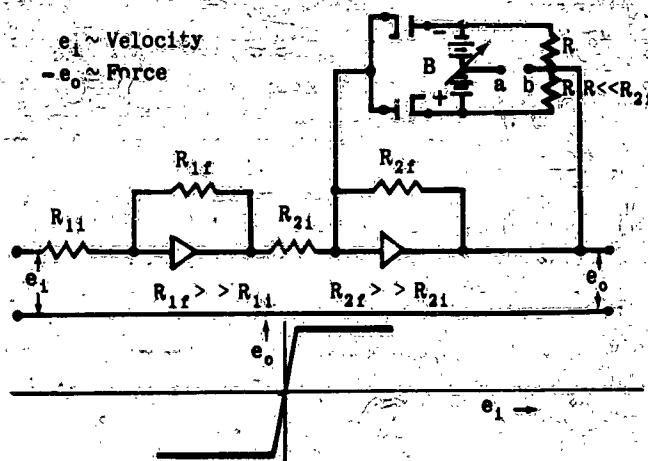


Figure VIII-15. Simulation of Coulomb Friction with Diodes

results in the conduction or lowered resistance of one of the diodes. Since this diode is in series with R, and the combination in parallel with the feedback resistor, the effect of a conducting diode is to lower the operational amplifier gain. This reduction in gain prevents the rise in output voltage which would normally be associated with a given rise in input voltage. As a result of the changing diode resistance and reduction of amplifier gain, the output voltage remains essentially constant. High gain amplifiers are used in this circuit for the same reason as with the polarized relay circuit.

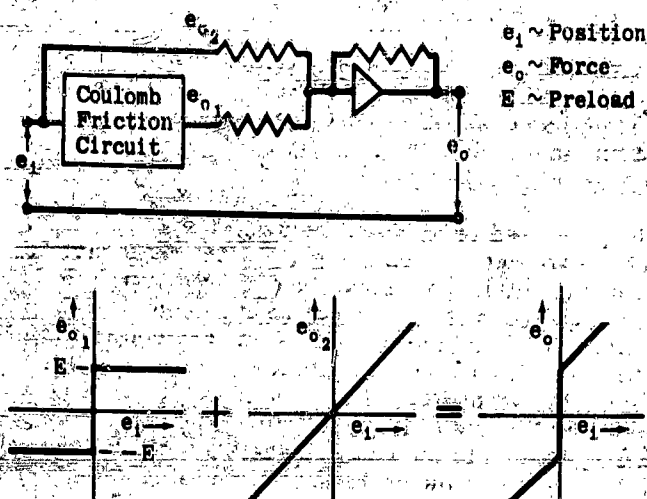


Figure VIII-16. Simulation of Spring Preload

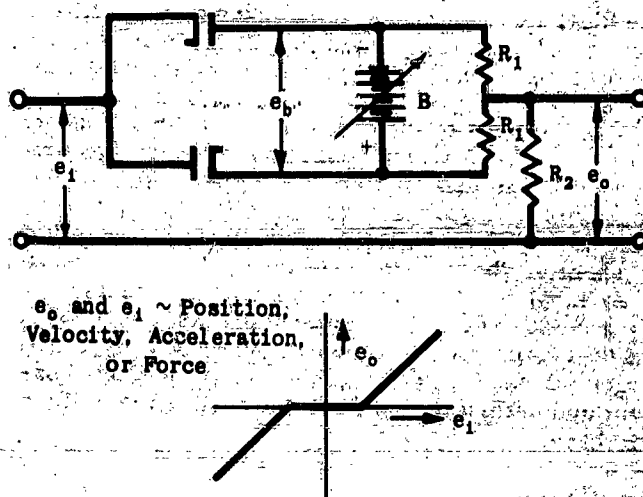


Figure VIII-17. Simulation of Threshold

(c) SPRING PRELOAD

To simulate spring preload, the circuits described for coulomb friction may be combined with one operational amplifier connected as an adder as shown in figure VIII-16.

(d) THRESHOLD

The simulation of threshold may be accomplished through a circuit similar to that of figure VIII-17. As shown by this figure, the diodes are non-conducting (very high resistance) when the input voltage is zero. For input voltages exceeding the bias voltage, one of the diodes conducts (lowered resistance) and the output voltage becomes $[R_2 / (R_1 + R_2 + R_0)] (e_1 - e_0)$. In this circuit, R_1 must be small compared to R_2 . Since the output of the diode circuit will interact with the input of the following operational amplifier, the threshold circuit must be connected to the circuit with which it is to operate when calibrating.

(e) LIMITING

A simple method of simulating limiting is shown in figure VIII-18. When R is small relative to the input impedance of the following stage, but still large compared to the resistance of a conducting diode, the output voltage will be that of the input until one of the diodes conducts and prevents a further increase in the output voltage.

(f) HYSTERESIS

The circuit of figure VIII-19 may be used to simulate hysteresis. As an aid in understanding the operation of this circuit, it may be recalled that the grid voltage of an operational amplifier is effectively zero. As may be seen by the figure, any positive input voltage large enough to cause one of the diodes to conduct will charge the condenser, C_1 , to a value equal to the input voltage minus one-half the bias battery voltage. As the input voltage begins to decrease, C_1 will maintain the voltage at point (A) at a constant value until the input voltage

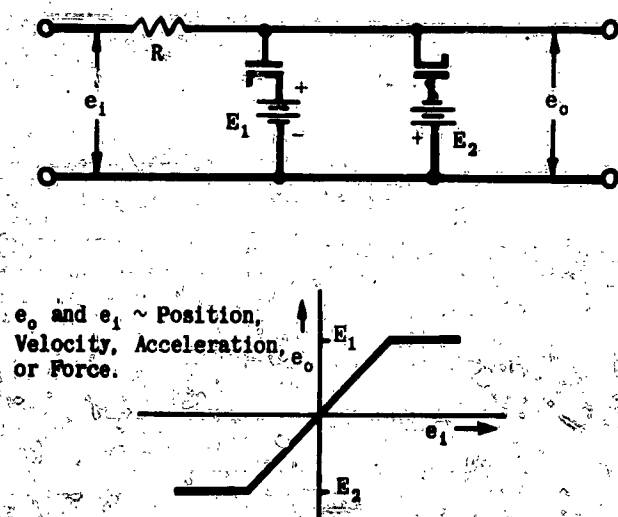


Figure VIII-18. Simulation of Limiting

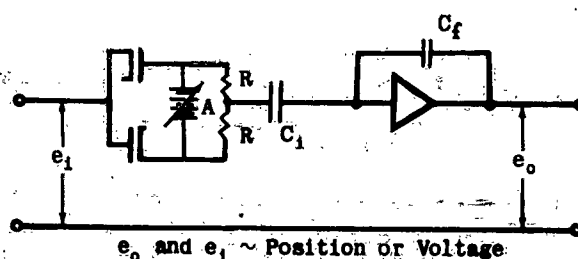


Figure VIII-19. Simulation of Hysteresis

has been reduced sufficiently to cause the other diode to conduct. When this second diode conducts, C_1 has a voltage equal to the input voltage plus the bias voltage. The feedback impedance of the operational amplifier is made a capacitive reactance to prevent differentiation due to the fact that the input impedance is a capacitive reactance.

(g) CONTINUOUS NON-LINEARITIES

Up to this point, the diodes have been used to simulate discontinuous non-linearities. In figure VIII-20 the use of diodes to approximate a smooth curve by straight lines is illustrated. This circuit is similar to figure VIII-15. The difference in this case lies in the resistances placed in series with each diode. These resistors permit the changing of the effective feedback impedance in steps. The diodes serve to switch in each resistor when the output voltage reaches a certain value.

The circuits described above are not necessarily the best ones for the particular non-linearities simulated. They are circuits which have been used; and serve to show how a few auxiliary components can be used to represent the majority of non-linearities.

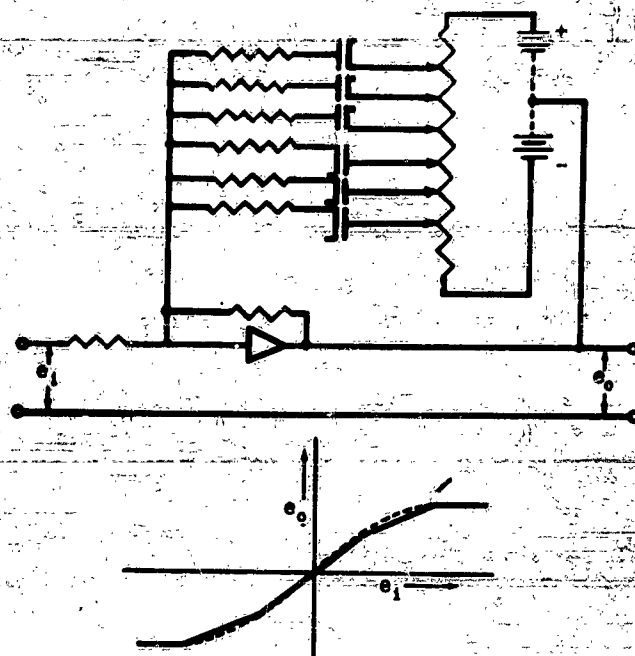


Figure VIII-20. Simulation of Continuous Non-Linearities

BIBLIOGRAPHY

The following bibliography is included for reference. The list is in no sense complete, but contains the major source material for this chapter. Many of the references, themselves, contain much more complete and detailed bibliographies.

1. 'Analog Computers for Servo Problems,' by D. McDonald; The Review of Scientific Instruments, Vol. 21, No. 2, February, 1950.
2. 'Linear Electronic Analog Computer Design,' by C. A. Meneley and C. D. Morrill; Goodyear Aircraft Corp., Akron, Ohio.
3. 'Computer Applications to Pilotless Aircraft Studies,' Report GER-2572 Goodyear Aircraft Corp., Akron, Ohio, 10 April 1951.
4. 'Electronic Analog Computers,' by G. A. Korn and T. M. Korn, McGraw-Hill, New York, 1952

APPENDIX

SECTION A — TABLES, CHARTS AND GRAPHS

The pages immediately following include certain charts that are useful in design. Some of them have been discussed in the text. Others not discussed are included because of their general interest.

PARAMETER		DEFINITION IN TERMS OF EQUATION COEFFICIENTS	EQUIVALENT EXPRESSIONS
DAMPING RATIO	ζ	$\frac{a_1}{2\sqrt{a_0 a_2}}$	$\frac{1}{\omega_n \tau}, \frac{T_n}{2\pi \tau}, \frac{\tau_1 + \tau_2}{2\sqrt{\tau_1 \tau_2}}, \frac{\nu + 1}{2\sqrt{\nu}}$
UNDAMPED ANGULAR NATURAL FREQUENCY	ω_n	$\sqrt{\frac{a_0}{a_2}}$	$\frac{\omega}{\sqrt{1-\zeta^2}}, 2\pi n_n, \frac{2\pi}{T_n}, \frac{1}{\zeta \tau}, \frac{1}{\sqrt{\tau_1 \tau_2}}$
UNDAMPED NATURAL FREQUENCY	n_n	$\frac{1}{2\pi\sqrt{\frac{a_2}{a_0}}}$	$\frac{\omega_n}{2\pi}, \frac{\omega}{2\pi\sqrt{1-\zeta^2}}, \frac{1}{T_n}, \frac{1}{2\pi\zeta\tau}$
UNDAMPED NATURAL PERIOD	T_n	$2\pi\sqrt{\frac{a_2}{a_0}}$	$\frac{2\pi}{\omega_n}, T\sqrt{1-\zeta^2}, 2\pi\zeta\tau, \frac{1}{n_n}$
ANGULAR NATURAL FREQUENCY	ω	$\sqrt{\frac{a_0}{a_2} - \left(\frac{a_1}{2a_2}\right)^2}$	$\omega_n\sqrt{1-\zeta^2}, 2\pi n, \frac{2\pi}{T}, \frac{\sqrt{1-\zeta^2}}{\zeta\tau}, \frac{2\pi\sqrt{1-\zeta^2}}{T_n}$
NATURAL FREQUENCY	n	$\frac{1}{2\pi}\sqrt{\frac{a_0}{a_2} - \left(\frac{a_1}{2a_2}\right)^2}$	$\frac{\omega_n\sqrt{1-\zeta^2}}{2\pi}, \frac{\omega}{2\pi}, \frac{1}{T}, \frac{\sqrt{1-\zeta^2}}{2\pi\zeta\tau}$
NATURAL PERIOD	T	$\frac{2\pi}{\sqrt{\frac{a_0}{a_2} - \left(\frac{a_1}{2a_2}\right)^2}}$	$\frac{2\pi}{\omega}, \frac{2\pi}{\omega_n\sqrt{1-\zeta^2}}, \frac{2\pi\zeta\tau}{\sqrt{1-\zeta^2}}, \frac{1}{n}, \frac{T_n}{\sqrt{1-\zeta^2}}$
CRITICAL TIME CONSTANT	τ	$\frac{2a_2}{a_1}$	$\frac{1}{\zeta\omega_n}, \frac{T_n}{2\pi\zeta}, \frac{T\sqrt{1-\zeta^2}}{2\pi\zeta}, \frac{2\tau_1\tau_2}{\tau_1 + \tau_2}, \frac{2\tau_1}{\nu + 1}, \frac{\tau_t}{2\zeta^2}$
LARGE TIME CONSTANT	τ_1	$\frac{1}{\frac{a_1}{2a_2} - \sqrt{\left(\frac{a_1}{2a_2}\right)^2 - \frac{a_0}{a_2}}}$	$\frac{1}{\omega_n[\zeta\sqrt{\zeta^2-1}]}, \frac{\sqrt{\nu}}{\omega_n}, \frac{T_n\sqrt{\nu}}{2\pi}, \sqrt{\nu}\zeta\tau, \nu\tau_2, \frac{\nu\tau_t}{\nu+1}$
SMALL TIME CONSTANT	τ_2	$\frac{1}{\frac{a_1}{2a_2} + \sqrt{\left(\frac{a_1}{2a_2}\right)^2 - \frac{a_0}{a_2}}}$	$\frac{1}{\omega_n[\zeta + \sqrt{\zeta^2-1}]}, \frac{1}{\sqrt{\nu}\omega_n}, \frac{T_n}{2\pi\sqrt{\nu}}, \frac{\zeta\tau}{\sqrt{\nu}}, \frac{\tau_1}{\nu}$
OVER-CRITICAL TIME CONSTANT	τ_t	$\frac{a_1}{a_0}$	$\frac{2\zeta}{\omega_n}, \tau_1 + \tau_2, \frac{\nu+1}{\nu}\tau_1, 2\zeta^2\tau, \frac{2\zeta}{\sqrt{\nu}}\tau$
TIME PARAMETER RATIO ($\zeta > 1$)	ν	$\frac{\frac{a_1}{2a_2} + \sqrt{\left(\frac{a_1}{2a_2}\right)^2 - \frac{a_0}{a_2}}}{\frac{a_1}{2a_2} - \sqrt{\left(\frac{a_1}{2a_2}\right)^2 - \frac{a_0}{a_2}}}$	$\frac{\zeta + \sqrt{\zeta^2-1}}{\zeta - \sqrt{\zeta^2-1}}, \frac{\tau_1}{\tau_2}$

Table A-1. Relationship Among System Parameters and Equation Coefficients for Second Order System
 $a_2[(d^2x)/(dt^2)] + a_1[(dx)/(dt)] + a_0x = 0$

Table A-2. Determination of Equation Coefficients for Second Order System From Response Curves

TRANSFER FUNCTION	TYPE OF RESPONSE	RESPONSE CURVE	EQUATION PARAMETERS USED	METHOD USED TO FIND EQUATION PARAMETERS	EQUATION COEFFICIENTS IN TERMS OF ω_n AND EQUATION PARAMETERS
$\frac{s}{s^2 + 2\zeta\omega_n s + \omega_n^2}$	Oscillatory $0 < \zeta < 0.5$		ζ, τ	Measure τ Form Ratios $\frac{x_1}{x_0}, \frac{x_2}{x_0}, \frac{x_3}{x_0}, \frac{x_4}{x_0}, \frac{x_5}{x_0}, \frac{x_6}{x_0}, \frac{x_7}{x_0}, \frac{x_8}{x_0}, \frac{x_9}{x_0}, \frac{x_{10}}{x_0}, \frac{x_{11}}{x_0}, \frac{x_{12}}{x_0}, \frac{x_{13}}{x_0}, \frac{x_{14}}{x_0}, \frac{x_{15}}{x_0}, \frac{x_{16}}{x_0}, \frac{x_{17}}{x_0}, \frac{x_{18}}{x_0}, \frac{x_{19}}{x_0}, \frac{x_{20}}{x_0}, \frac{x_{21}}{x_0}, \frac{x_{22}}{x_0}, \frac{x_{23}}{x_0}, \frac{x_{24}}{x_0}, \frac{x_{25}}{x_0}, \frac{x_{26}}{x_0}, \frac{x_{27}}{x_0}, \frac{x_{28}}{x_0}, \frac{x_{29}}{x_0}, \frac{x_{30}}{x_0}, \frac{x_{31}}{x_0}, \frac{x_{32}}{x_0}, \frac{x_{33}}{x_0}, \frac{x_{34}}{x_0}, \frac{x_{35}}{x_0}, \frac{x_{36}}{x_0}, \frac{x_{37}}{x_0}, \frac{x_{38}}{x_0}, \frac{x_{39}}{x_0}, \frac{x_{40}}{x_0}, \frac{x_{41}}{x_0}, \frac{x_{42}}{x_0}, \frac{x_{43}}{x_0}, \frac{x_{44}}{x_0}, \frac{x_{45}}{x_0}, \frac{x_{46}}{x_0}, \frac{x_{47}}{x_0}, \frac{x_{48}}{x_0}, \frac{x_{49}}{x_0}, \frac{x_{50}}{x_0}, \frac{x_{51}}{x_0}, \frac{x_{52}}{x_0}, \frac{x_{53}}{x_0}, \frac{x_{54}}{x_0}, \frac{x_{55}}{x_0}, \frac{x_{56}}{x_0}, \frac{x_{57}}{x_0}, \frac{x_{58}}{x_0}, \frac{x_{59}}{x_0}, \frac{x_{60}}{x_0}, \frac{x_{61}}{x_0}, \frac{x_{62}}{x_0}, \frac{x_{63}}{x_0}, \frac{x_{64}}{x_0}, \frac{x_{65}}{x_0}, \frac{x_{66}}{x_0}, \frac{x_{67}}{x_0}, \frac{x_{68}}{x_0}, \frac{x_{69}}{x_0}, \frac{x_{70}}{x_0}, \frac{x_{71}}{x_0}, \frac{x_{72}}{x_0}, \frac{x_{73}}{x_0}, \frac{x_{74}}{x_0}, \frac{x_{75}}{x_0}, \frac{x_{76}}{x_0}, \frac{x_{77}}{x_0}, \frac{x_{78}}{x_0}, \frac{x_{79}}{x_0}, \frac{x_{80}}{x_0}, \frac{x_{81}}{x_0}, \frac{x_{82}}{x_0}, \frac{x_{83}}{x_0}, \frac{x_{84}}{x_0}, \frac{x_{85}}{x_0}, \frac{x_{86}}{x_0}, \frac{x_{87}}{x_0}, \frac{x_{88}}{x_0}, \frac{x_{89}}{x_0}, \frac{x_{90}}{x_0}, \frac{x_{91}}{x_0}, \frac{x_{92}}{x_0}, \frac{x_{93}}{x_0}, \frac{x_{94}}{x_0}, \frac{x_{95}}{x_0}, \frac{x_{96}}{x_0}, \frac{x_{97}}{x_0}, \frac{x_{98}}{x_0}, \frac{x_{99}}{x_0}, \frac{x_{100}}{x_0}$ Find ζ from Fig. A-4	$\frac{\omega_n}{2\zeta}$ $\frac{\omega_n^2}{1 - \zeta^2}$
$\frac{s}{s^2 + 2\zeta\omega_n s + \omega_n^2}$	Near critically Aperiodic $0.5 < \zeta < 2.0$		ζ, ω_n	Measure t_1, t_2, t_3 Form Ratios $\frac{t_1}{t_2}, \frac{t_1}{t_3}, \frac{t_2}{t_3}$ Find ζ, ω_n from Fig. A-5 Compute value of ω_n from $\zeta = \frac{\omega_n}{\omega_n}$	$\frac{\omega_n}{2\zeta}$ $\frac{\omega_n^2}{1 - \zeta^2}$
$\frac{s}{s^2 + 2\zeta\omega_n s + \omega_n^2}$	Critically Aperiodic $\zeta = 1.0$		τ	[Measure τ on Fig. A-5] (For $\zeta = 1.0, \tau = t_1 - t_2 - t_3 - t_4$)	$2\omega_n \tau$
$\frac{s}{s^2 + 2\zeta\omega_n s + \omega_n^2}$	Non Oscillatory $\zeta > 1.0$		ν, τ_1	Measure ν and τ_1 at t_1 & t_2 respectively. Measure x_0 & $x(\infty)$ Compute τ_1 from $\tau_1 = \frac{\ln x_1 - \ln x_2}{\ln x_1 - \ln x_2}$ Compute ν from $\nu = \frac{x(\infty) - x_0}{x(\infty)}$ Plot Response Curve on Semilog Paper. Extrapolate straight line portion of plot to $t = 0$.	$\frac{\omega_n}{2\zeta}$ $\frac{\omega_n^2}{1 - \zeta^2}$

	$F(s)$	$f(t)$	TIME RESPONSE
STEP FUNCTION POSITION	① $\frac{1}{s}$	1	
STEP FUNCTION VELOCITY	② $\frac{1}{s^2}$	t	
STEP FUNCTION ACCELERATION	③ $\frac{1}{s^3}$	$\frac{1}{2} t^2$	
FIRST ORDER LAG CONVERGING	④ $\frac{1}{\tau s + 1}$	$\frac{1}{\tau} e^{-\frac{t}{\tau}}$	
	⑤ $\frac{1}{s(\tau s + 1)}$	$1 - e^{-\frac{t}{\tau}}$	
FIRST ORDER LAG DIVERGING	⑥ $\frac{1}{-\tau s + 1}$	$\frac{1}{\tau} e^{\frac{t}{\tau}}$	
	⑦ $\frac{1}{s(-\tau s + 1)}$	$-1 + e^{\frac{t}{\tau}}$	
UNDAMPED SECOND ORDER $\zeta = 0$	⑧ $\frac{1}{\frac{s^2}{\omega_n^2} + 1}$ $\zeta = 0$	$\omega_n \sin \omega_n t$	
	⑨ $\frac{s}{\frac{s^2}{\omega_n^2} + 1}$ $\zeta = 0$	$\omega_n^2 \cos \omega_n t$	

Table A-3 (Sheet 1 of 2 Sheets). Time Responses
of Some Common Transient Modes

	$F(s)$	$f(t)$	TIME RESPONSE
CONVERGING SECOND ORDER $0 < \zeta \leq 1$	10 $\zeta < 1$ $\frac{1}{\frac{s^2}{\omega_n^2} + 2\frac{\zeta}{\omega_n}s + 1}$	$\frac{\omega_n}{\sqrt{1-\zeta^2}} e^{-\zeta\omega_n t} \sin \omega_n \sqrt{1-\zeta^2} t$	
	11 $\zeta < 1$ $\frac{s}{\frac{s^2}{\omega_n^2} + 2\frac{\zeta}{\omega_n}s + 1}$	$\frac{\omega_n^2}{\sqrt{1-\zeta^2}} e^{-\zeta\omega_n t} \sin(\omega_n \sqrt{1-\zeta^2} t + \tan^{-1} \frac{\sqrt{1-\zeta^2}}{-\zeta})$	
	12 $\zeta < 1$ $\frac{1}{s(\frac{s^2}{\omega_n^2} + 2\frac{\zeta}{\omega_n}s + 1)}$	$1 - \frac{e^{-\zeta\omega_n t}}{\sqrt{1-\zeta^2}} \sin(\omega_n \sqrt{1-\zeta^2} t + \tan^{-1} \frac{\sqrt{1-\zeta^2}}{-\zeta})$	
	13 $\zeta = 1$ $\frac{1}{\frac{s^2}{\omega_n^2} + 2\frac{\zeta}{\omega_n}s + 1}$	$\omega_n^2 t e^{-\omega_n t}$	
DIVERGING SECOND ORDER $\zeta < 0$ $\zeta' = \zeta $	14 $\zeta < 0$ $\frac{1}{\frac{s^2}{\omega_n^2} + 2\frac{\zeta}{\omega_n}s + 1}$	$\frac{\omega_n}{\sqrt{1-\zeta'^2}} e^{\zeta'\omega_n t} \sin \omega_n \sqrt{1-\zeta'^2} t$	
	15 $\zeta < 0$ $\frac{s}{\frac{s^2}{\omega_n^2} + 2\frac{\zeta}{\omega_n}s + 1}$	$\frac{\omega_n^2}{\sqrt{1-\zeta'^2}} e^{\zeta'\omega_n t} \sin(\omega_n \sqrt{1-\zeta'^2} t + \tan^{-1} \frac{\sqrt{1-\zeta'^2}}{\zeta'})$	
	16 $\zeta < 0$ $\frac{1}{s(\frac{s^2}{\omega_n^2} + 2\frac{\zeta}{\omega_n}s + 1)}$	$1 - \frac{e^{\zeta'\omega_n t}}{\sqrt{1-\zeta'^2}} \sin(\omega_n \sqrt{1-\zeta'^2} t + \tan^{-1} \frac{\sqrt{1-\zeta'^2}}{\zeta'})$	

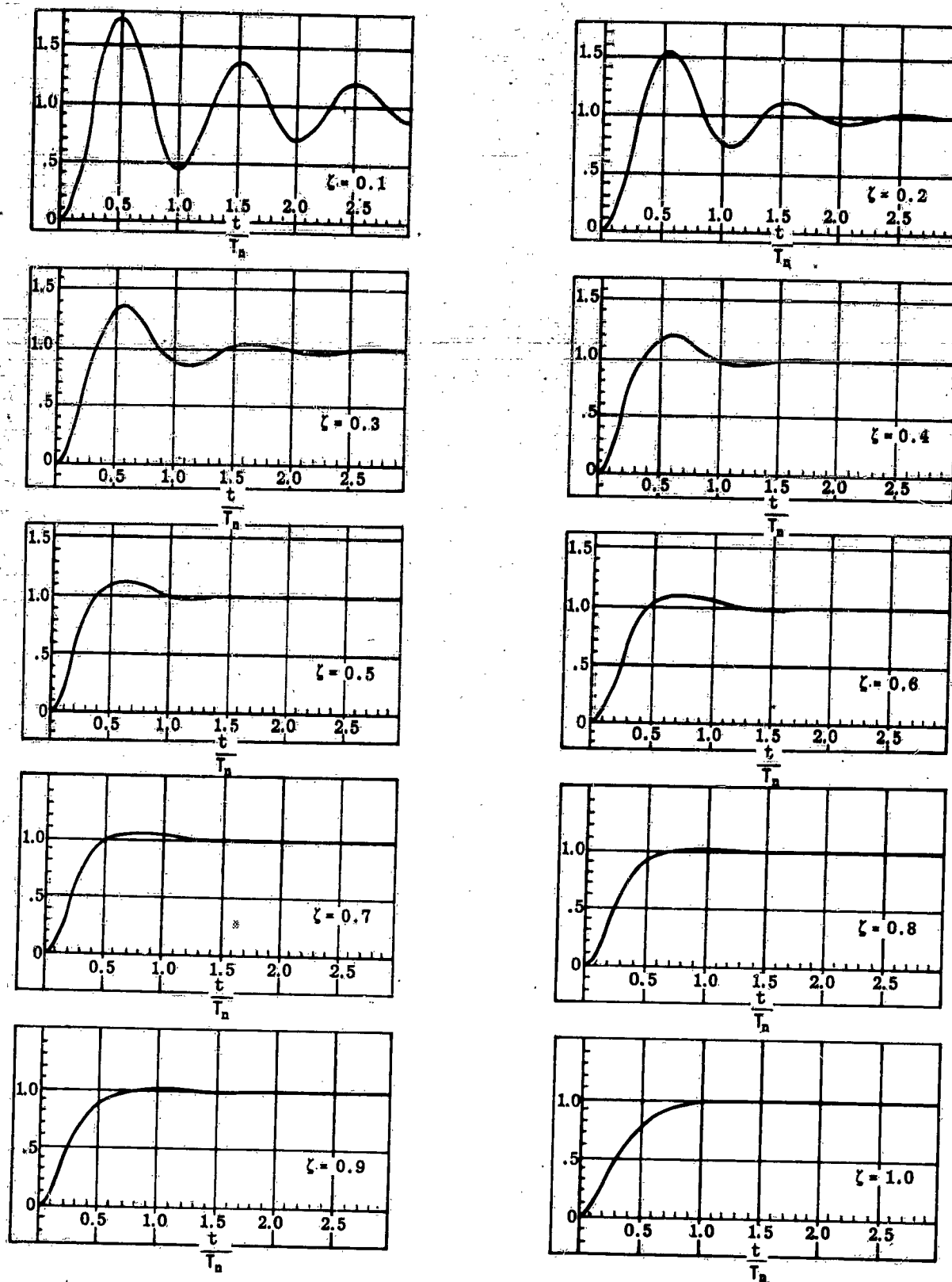
Table A-3 (Sheet 2 of 2 Sheets). Time Responses of Some Common Transient Modes

TRANSFER FUNCTION $Y(s)$	REPRESENTATION IN s -PLANE	REPRESENTATION IN COMPLEX $Y(s)$ PLANE	LOGARITHMIC REPRESENTATION
1. s			
2. s^2			
3. s^3			
4. $\frac{1}{s}$			
5. $\frac{1}{s^2}$			
6. $\frac{1}{s^3}$			
7. $\tau s + 1$			
8. $\frac{1}{\tau s + 1}$			

Table A-4 (Sheet 1 of 2 Sheets). Summary of the Forms of the Factors of the Transfer Function

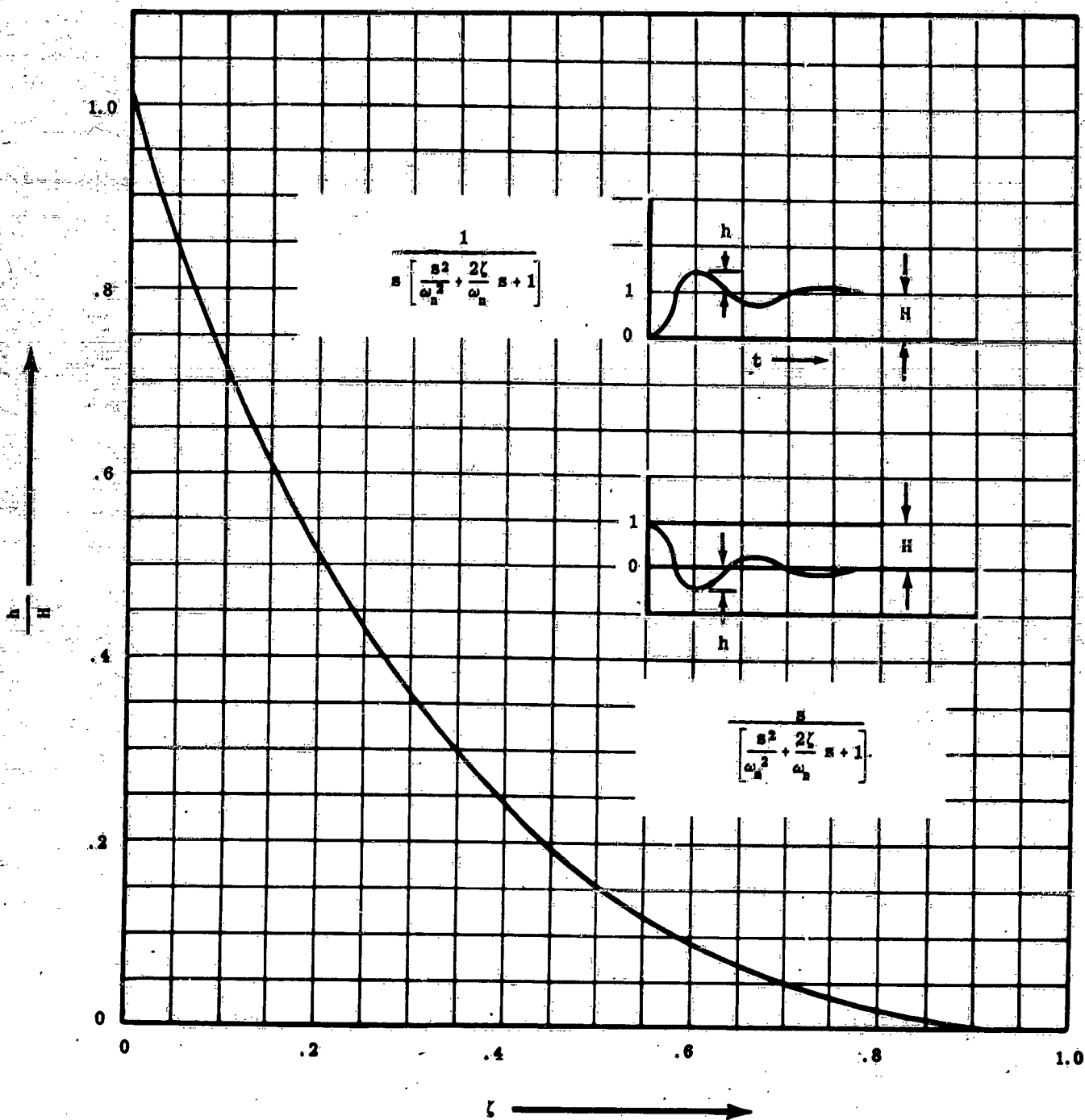
TRANSFER FUNCTION $Y(s)$	REPRESENTATION IN s -PLANE	REPRESENTATION IN COMPLEX $Y(s)$ PLANE	LOGARITHMIC REPRESENTATION
9. $\frac{-\tau s + 1}{\omega_n^2}$ (Non-minimum Phase)			
10. $\frac{1}{-\tau s + 1}$ (Non-minimum Phase)			
11. $\frac{s^2 + 2\zeta s + 1}{\omega_n^2}$			
12. $\frac{1}{s^2 + 2\zeta s + 1}$			
13. $\frac{s^2 - 2\zeta s + 1}{\omega_n^2}$ (Non-Minimum Phase)			
14. $\frac{1}{s^2 - 2\zeta s + 1}$ (Non-Minimum Phase)			

Table A-4 (Sheet 2 of 2 Sheets). Summary of the Forms
of the Factors of the Transfer Function



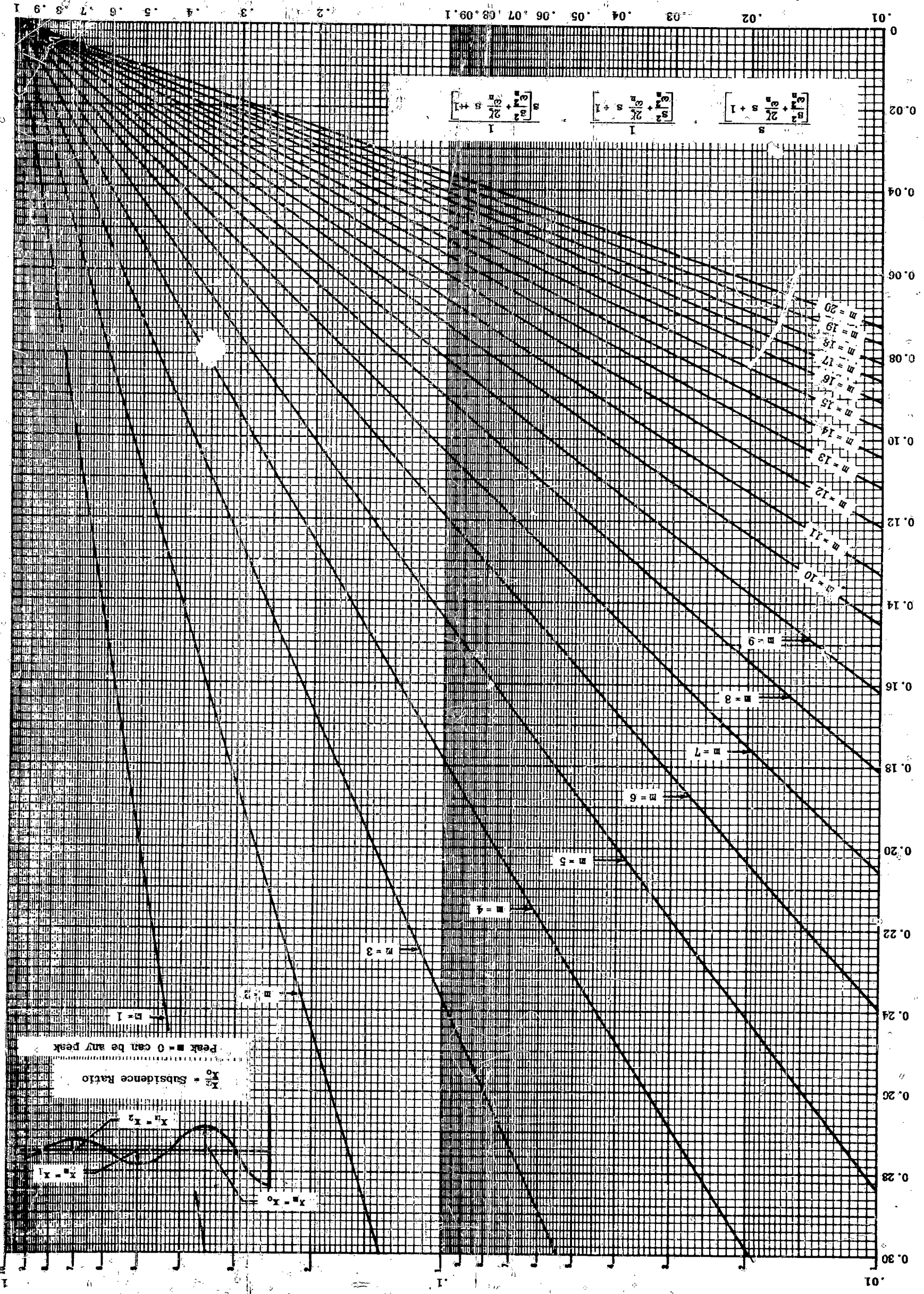
$$\frac{1}{s \left[\frac{s^2}{\omega_n^2} + \frac{2\zeta}{\omega_n} s + 1 \right]}$$

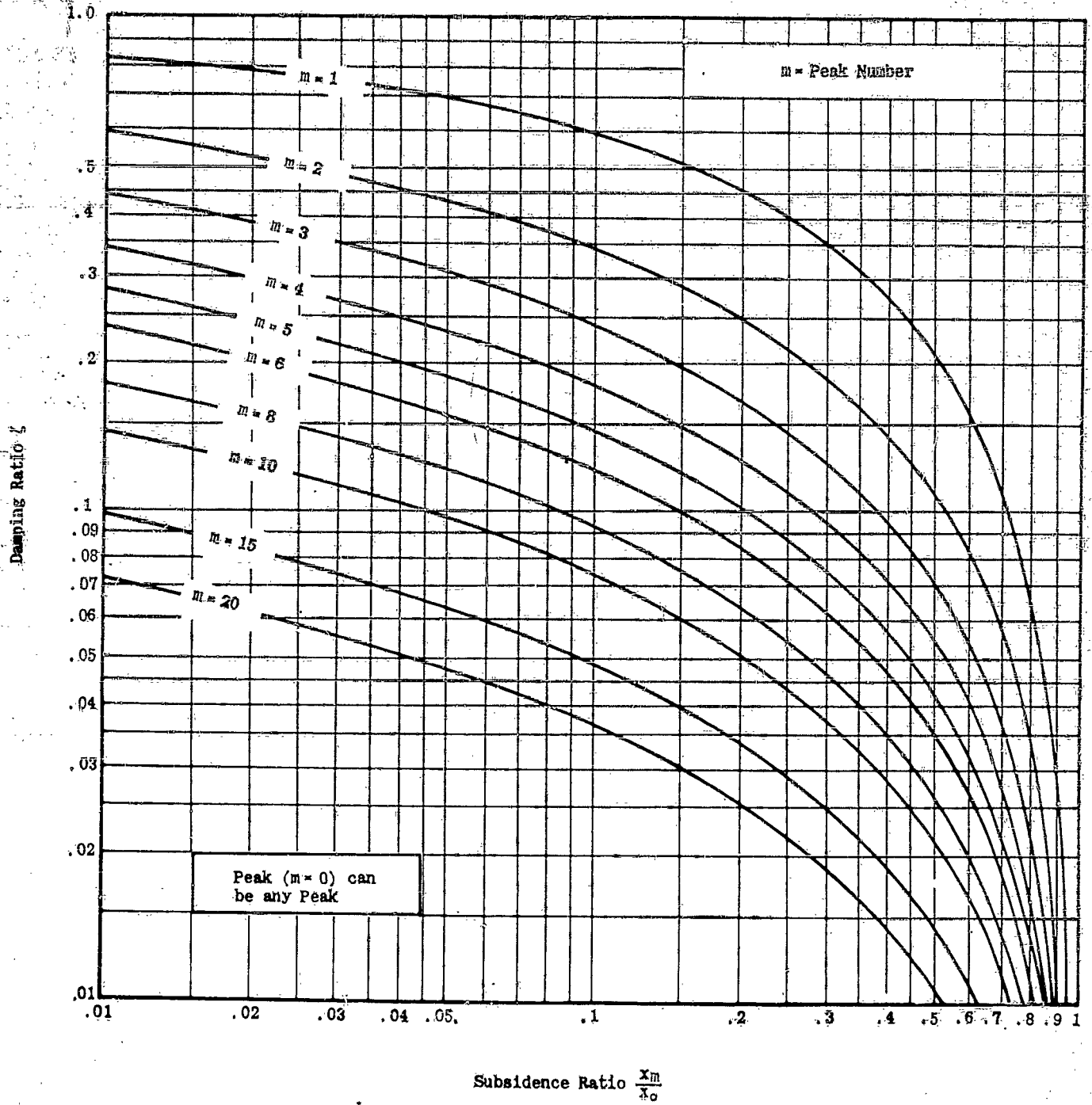
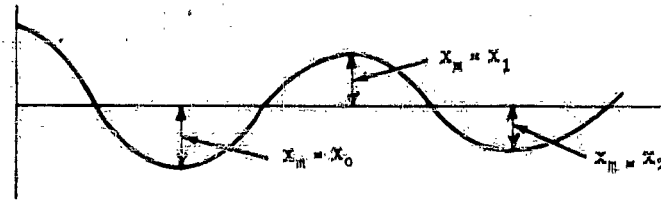
Figure A-1. Typical Response Curves of Second Order System to Step Function Disturbance When Damping Ratio is Less than Critical ($\zeta < 1$)



A-2. Determination of Equation Coefficients for Second Order Systems from Response Curves

Figure A-3. Damping Ratio of Oscillatory Transients as a Function of Subsidence Ratio for Second Order Systems





$$\frac{s}{\frac{s^2}{\omega_n^2} + \frac{2\zeta}{\omega_n} s + 1}, \frac{1}{\frac{s^2}{\omega_n^2} + \frac{2\zeta}{\omega_n} s + 1}, \frac{1}{s \left[\frac{s^2}{\omega_n^2} + \frac{2\zeta}{\omega_n} s + 1 \right]}$$

Figure A-4. Damping Ratio of Oscillatory Transients as a Function of Subsidence Ratio for Second Order Systems

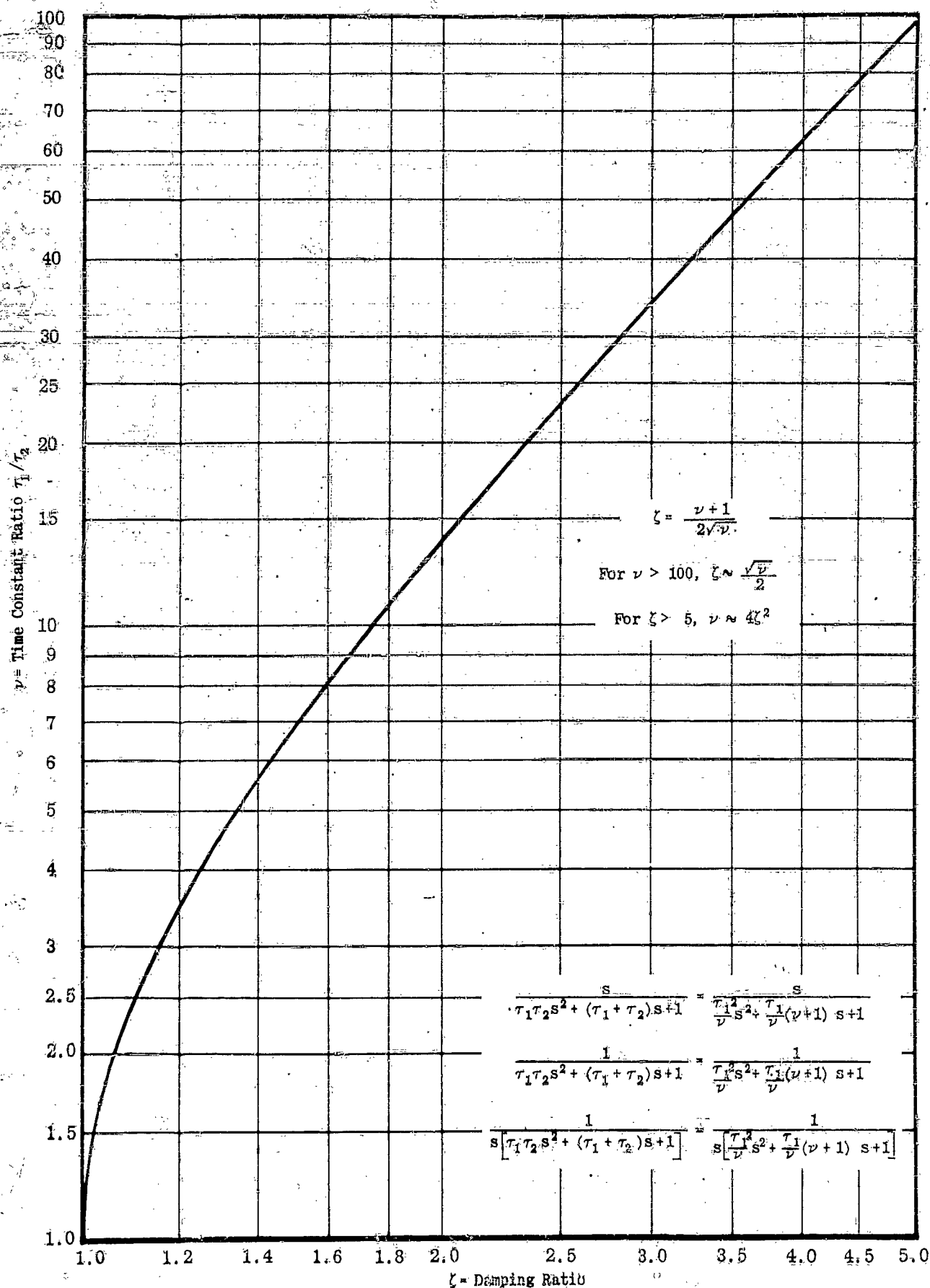
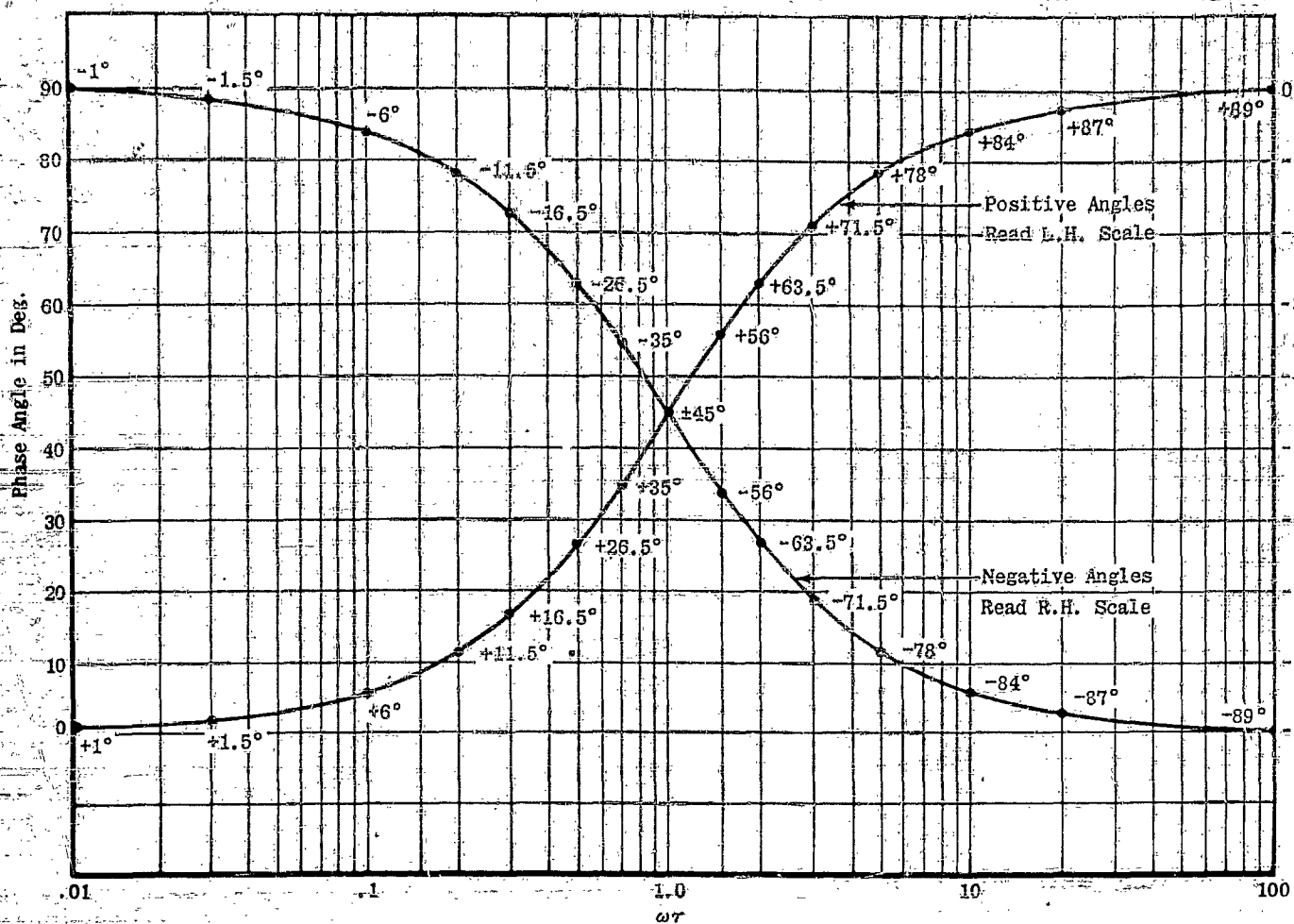


Figure A-6. Time Constant Ratio τ_1/τ_2 as a Function of Damping Ratio for Overdamped Second Order System



$$(\pm rs + 1)^{\pm 1}$$

$$s = j\omega$$

$$\text{Phase Angle} = \tan^{-1} \omega\tau$$

Figure A-7. Phase Angle for First Order System.

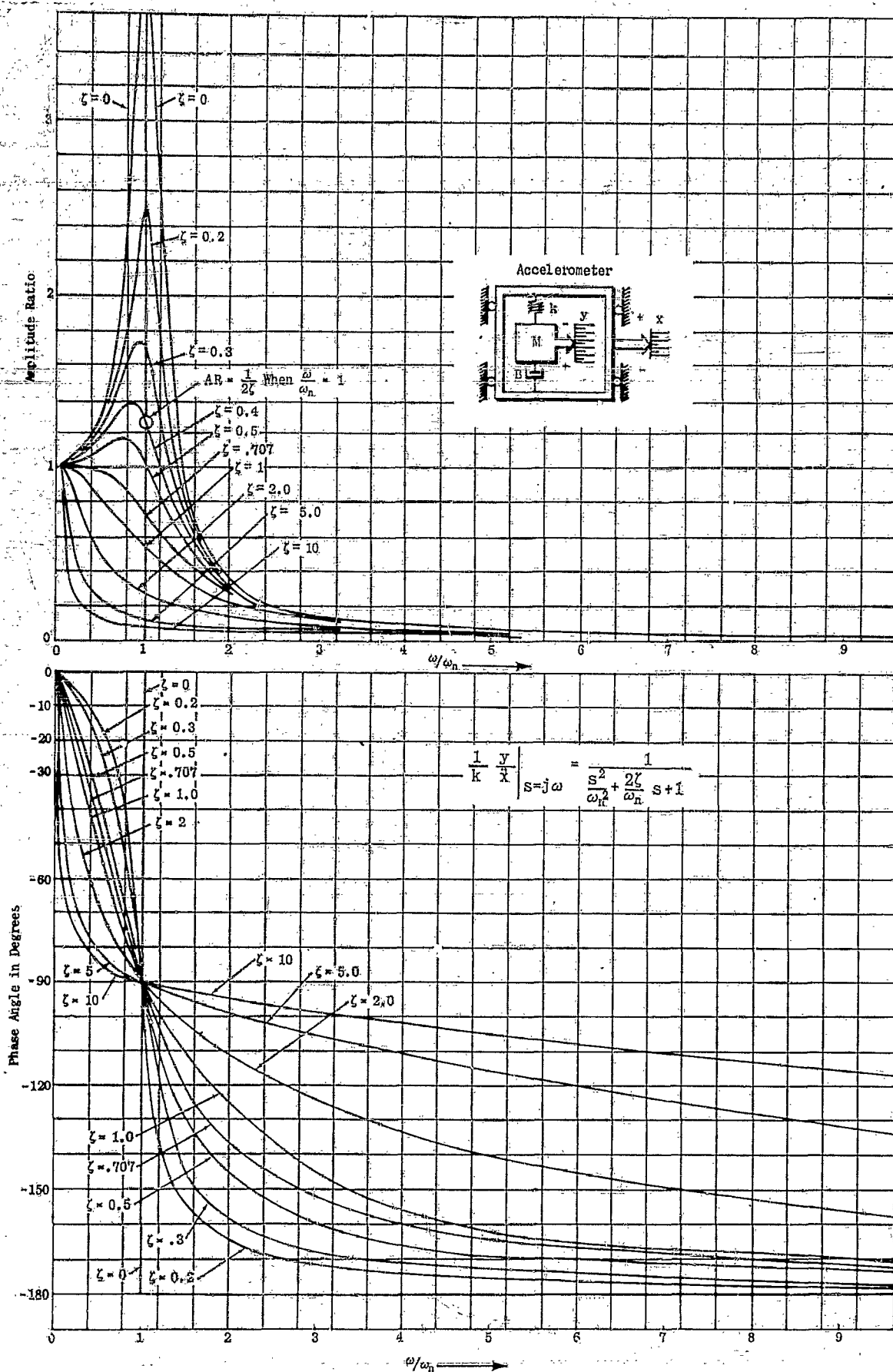


Figure A-8. Response of a Second Order Accelerometer to Sinusoidal Displacement Input.

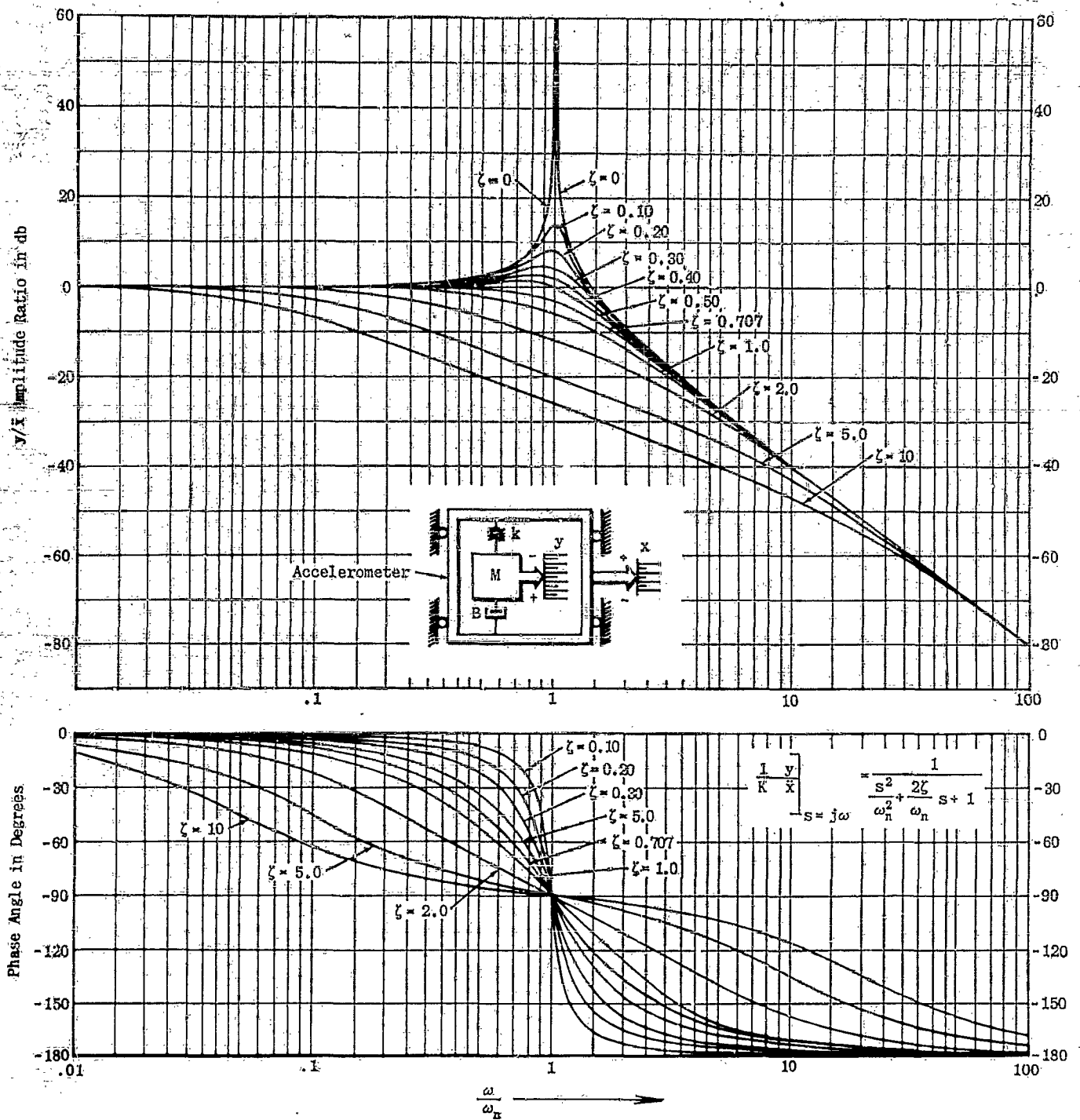
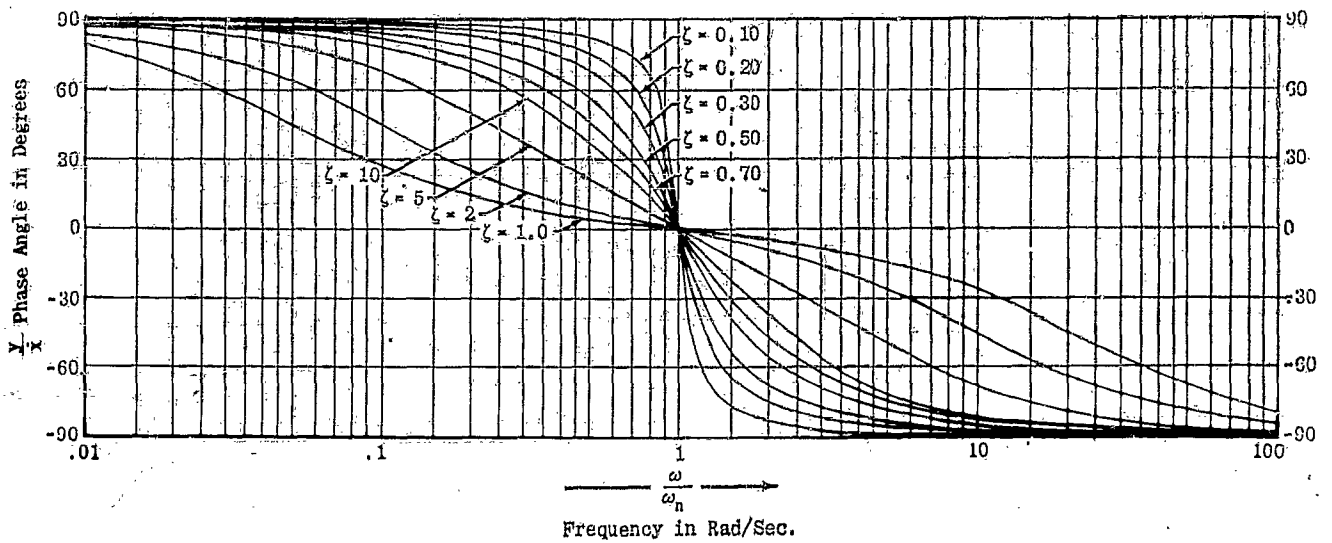
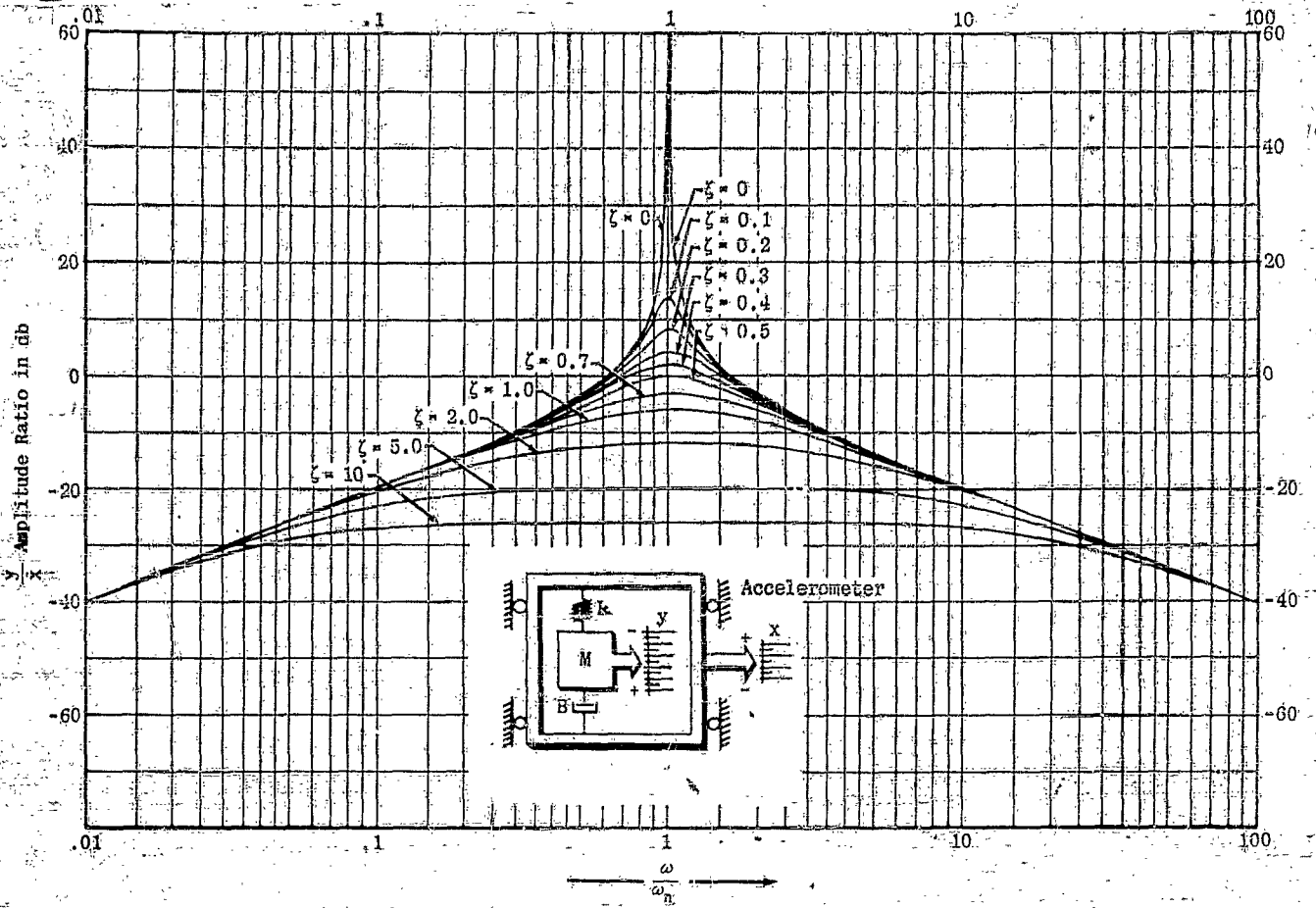


Figure A-9. Response of a Second Order Accelerometer to Sinusoidal Acceleration Input



$$\frac{1}{K} \frac{y}{x} = \frac{s}{\frac{s^2}{\omega_n^2} + \frac{2\zeta}{\omega_n} s + 1}$$

$s = j\omega$

Figure A-10. Response of a Second Order Accelerometer to Sinusoidal Velocity Input

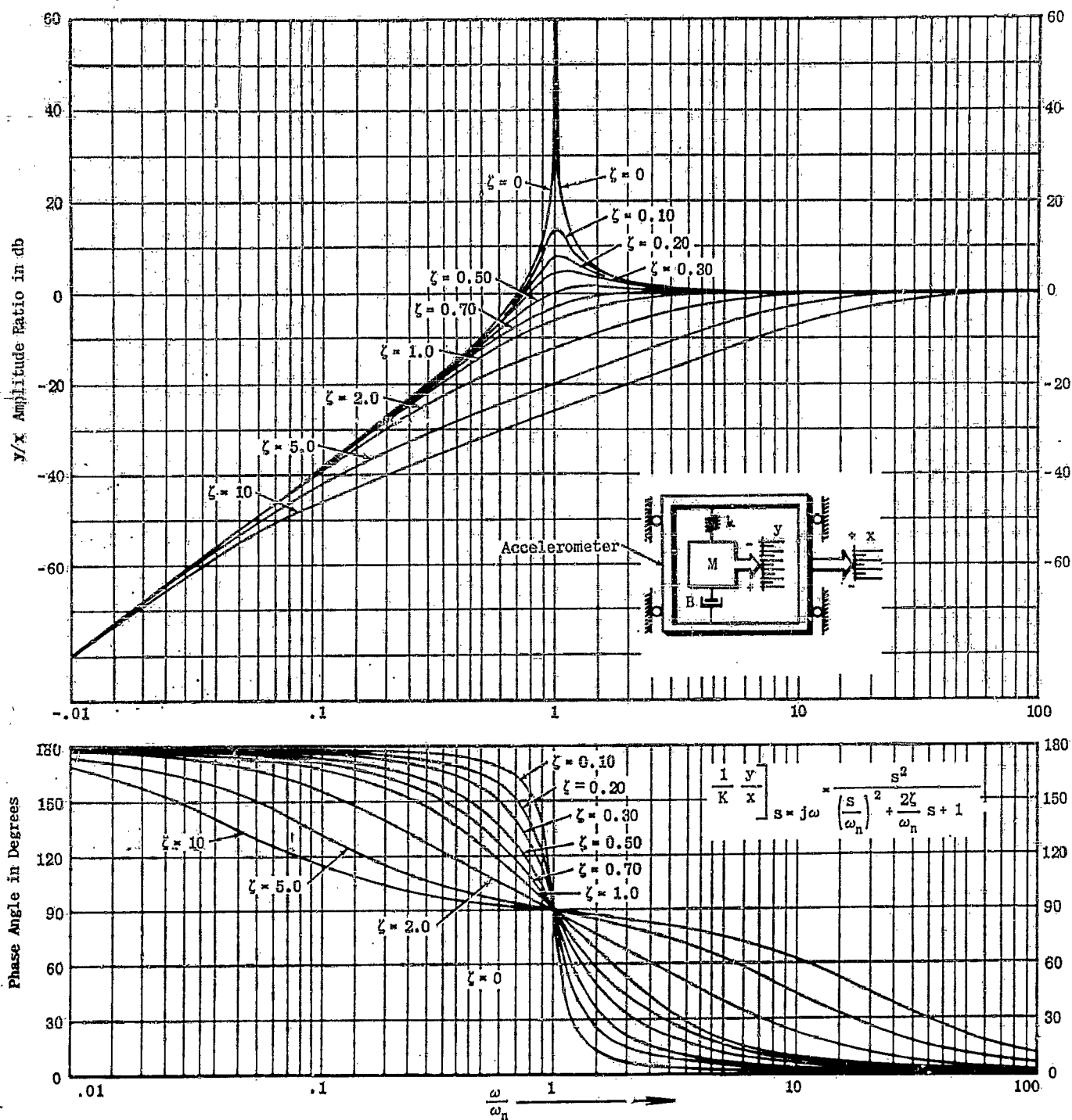


Figure A-11. Response of a Second Order Accelerometer to Sinusoidal Displacement Input

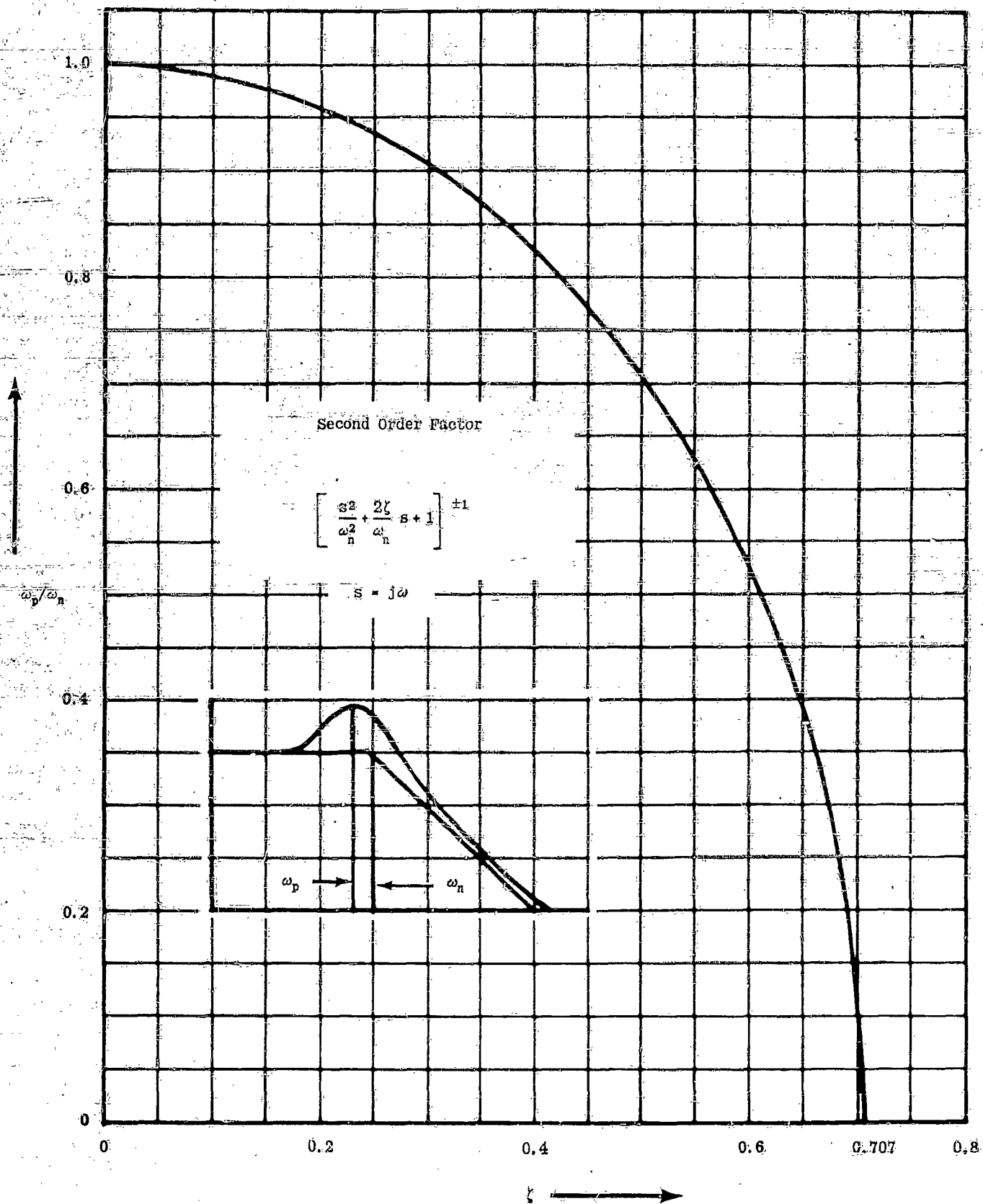
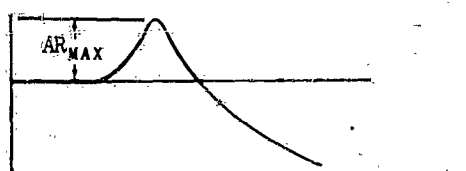
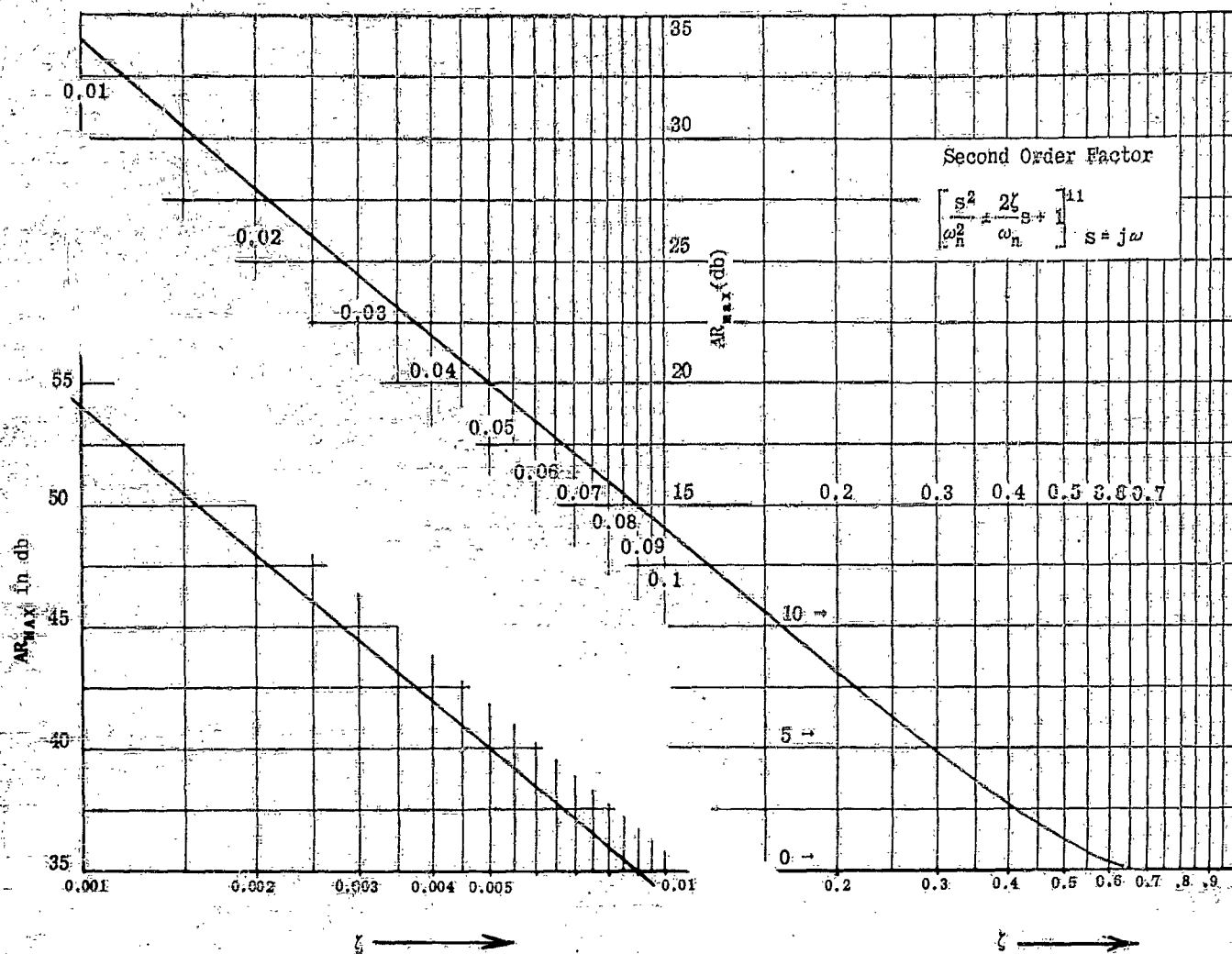
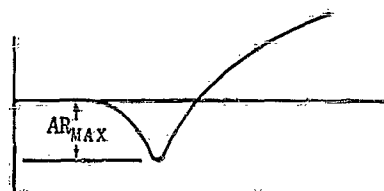


Figure A-12. Frequency at which AR_{MAX} Occurs



$$\left[\frac{s^2}{\omega_n^2} + \frac{2\zeta}{\omega_n} s + 1 \right]^{-1}$$



$$\left[\frac{s^2}{\omega_n^2} + \frac{2\zeta}{\omega_n} s + 1 \right]$$

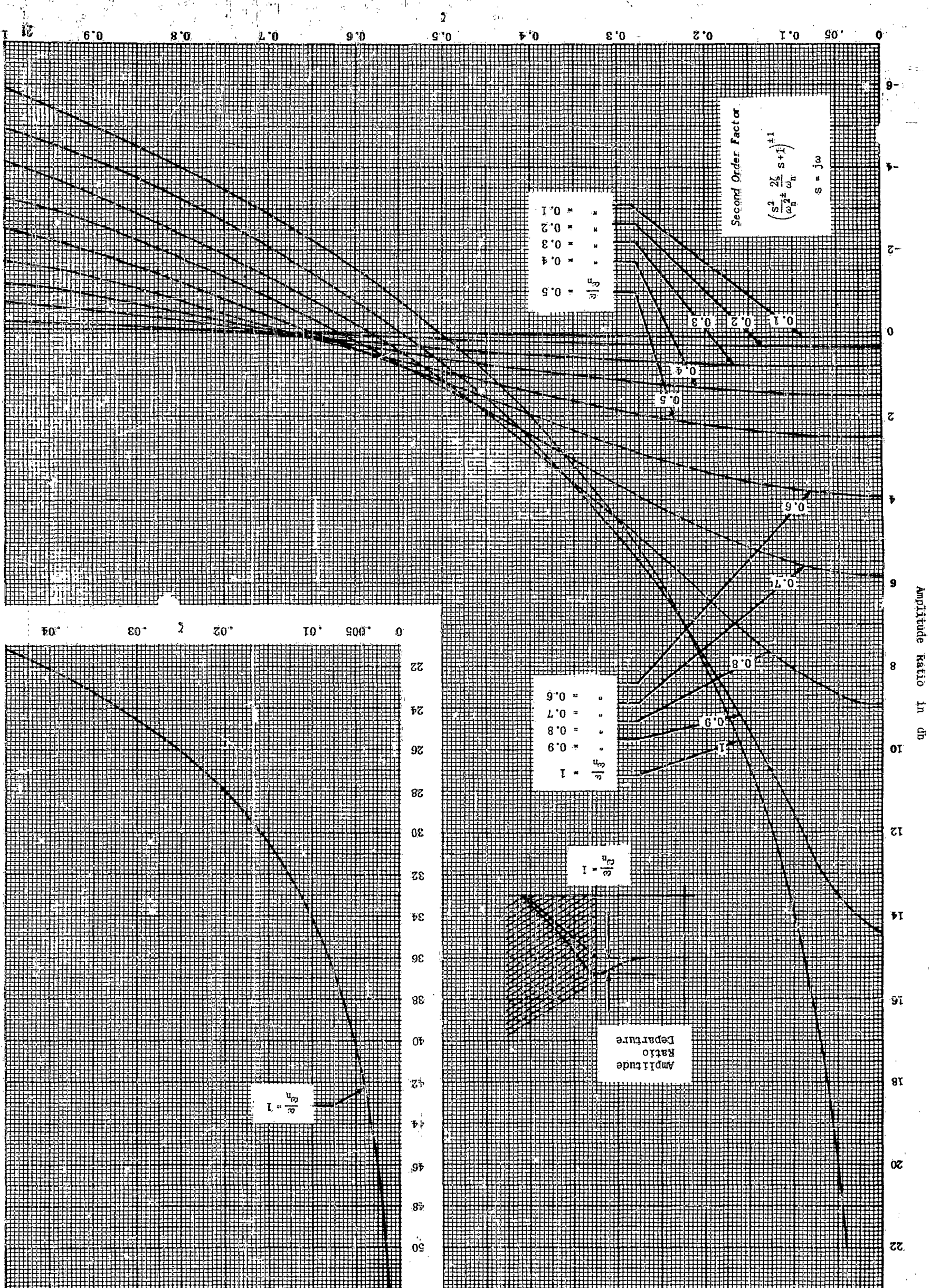
Second Order Factor

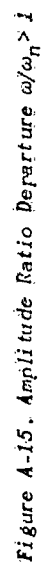
$$\left[\frac{s^2}{\omega_n^2} + \frac{2\zeta}{\omega_n} s + 1 \right]^{\pm 1}$$

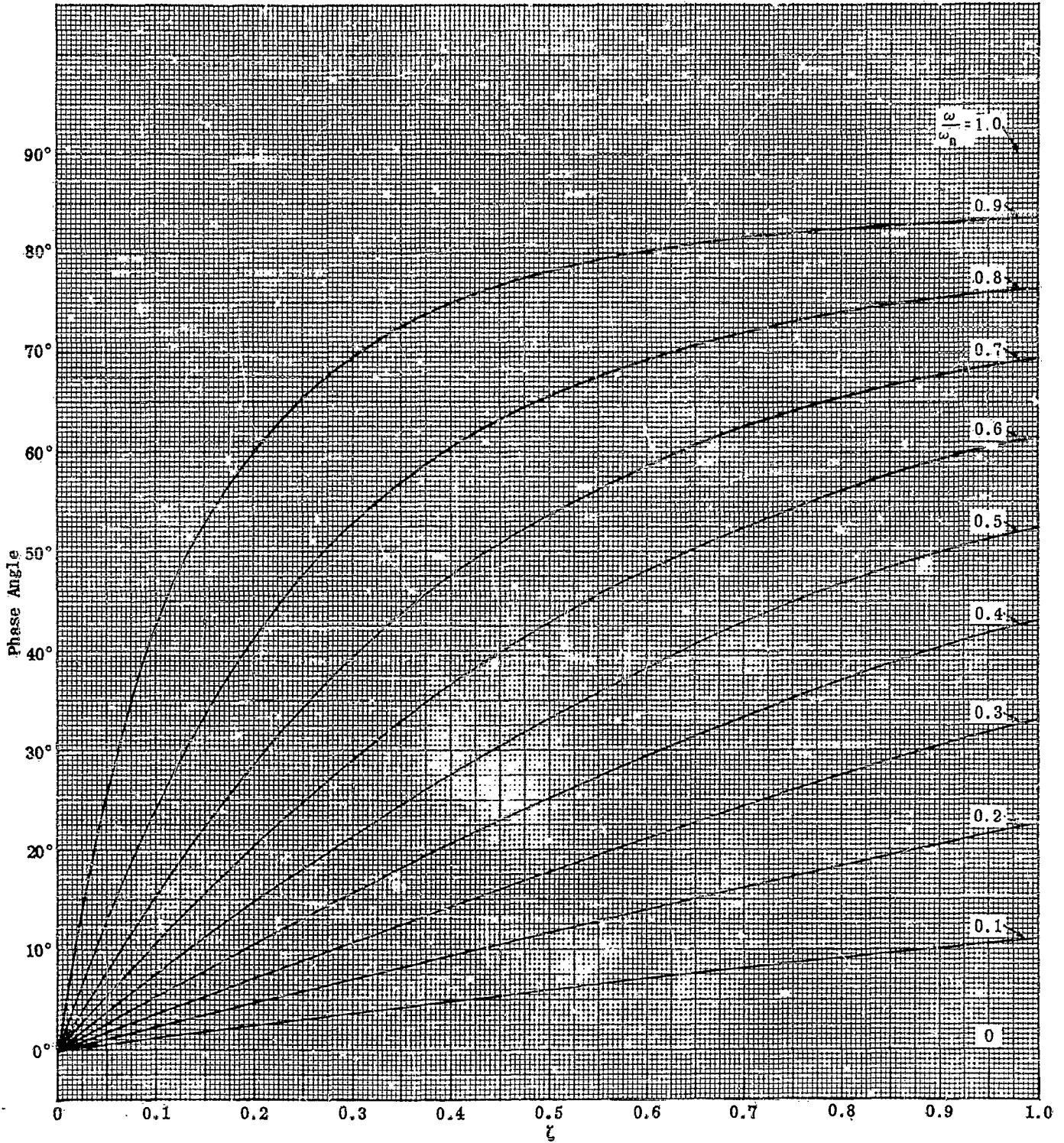
$$s = j\omega$$

Figure A-13. Plot of AR_{MAX}

Figure A-14. Amplitude Ratio Departure $\omega/\omega_n \leq 1$







Second Order Factor

$$\left[\frac{s^2 + \frac{2\zeta}{\omega_n} s + 1}{\omega_n^2} \right]^{+1}$$

$$s = j\omega$$

Figure A-16. Phase Angle $\omega/\omega_n < 1$

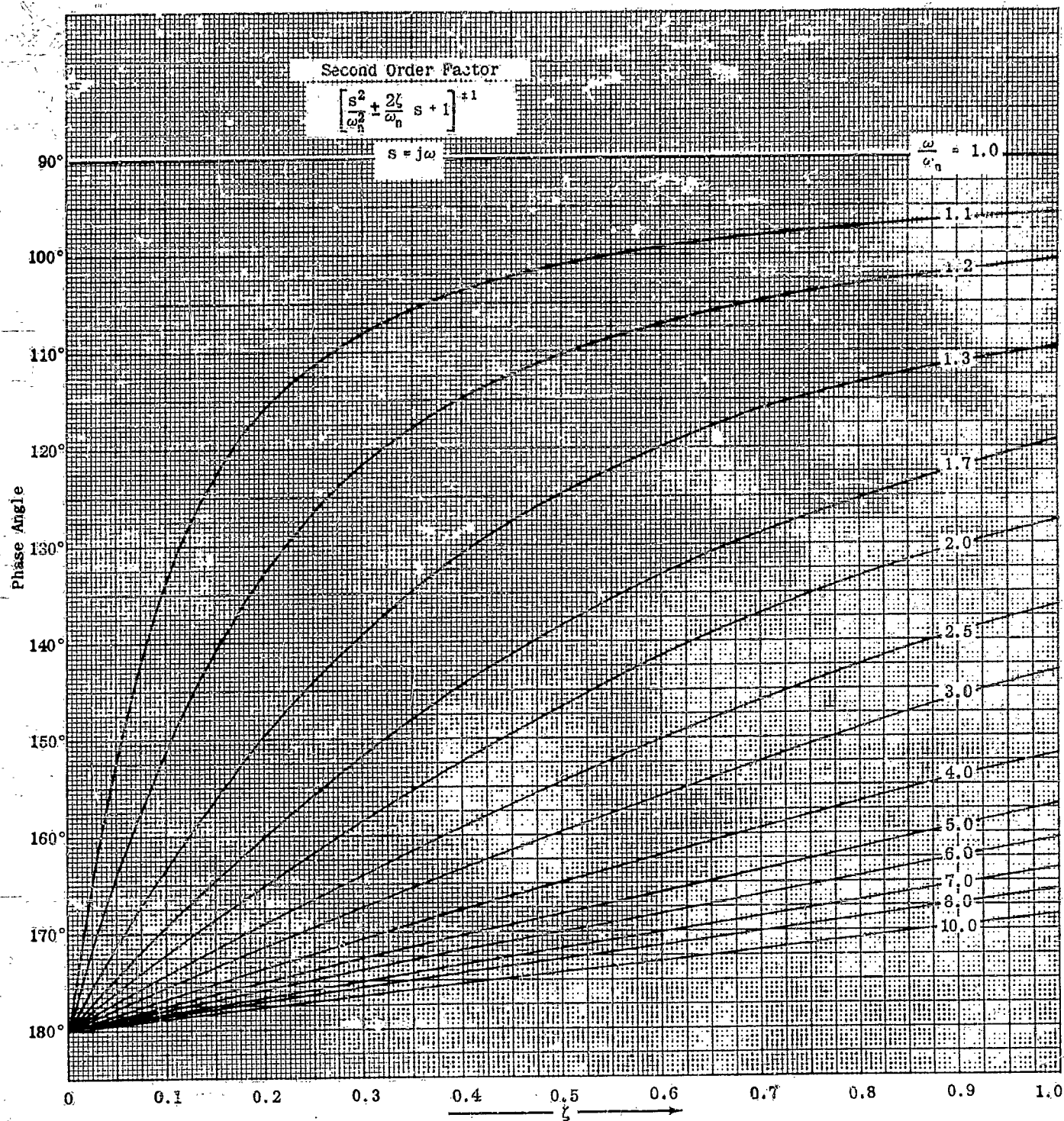


Figure A-17. Phase Angle $\omega/\omega_n > 1$

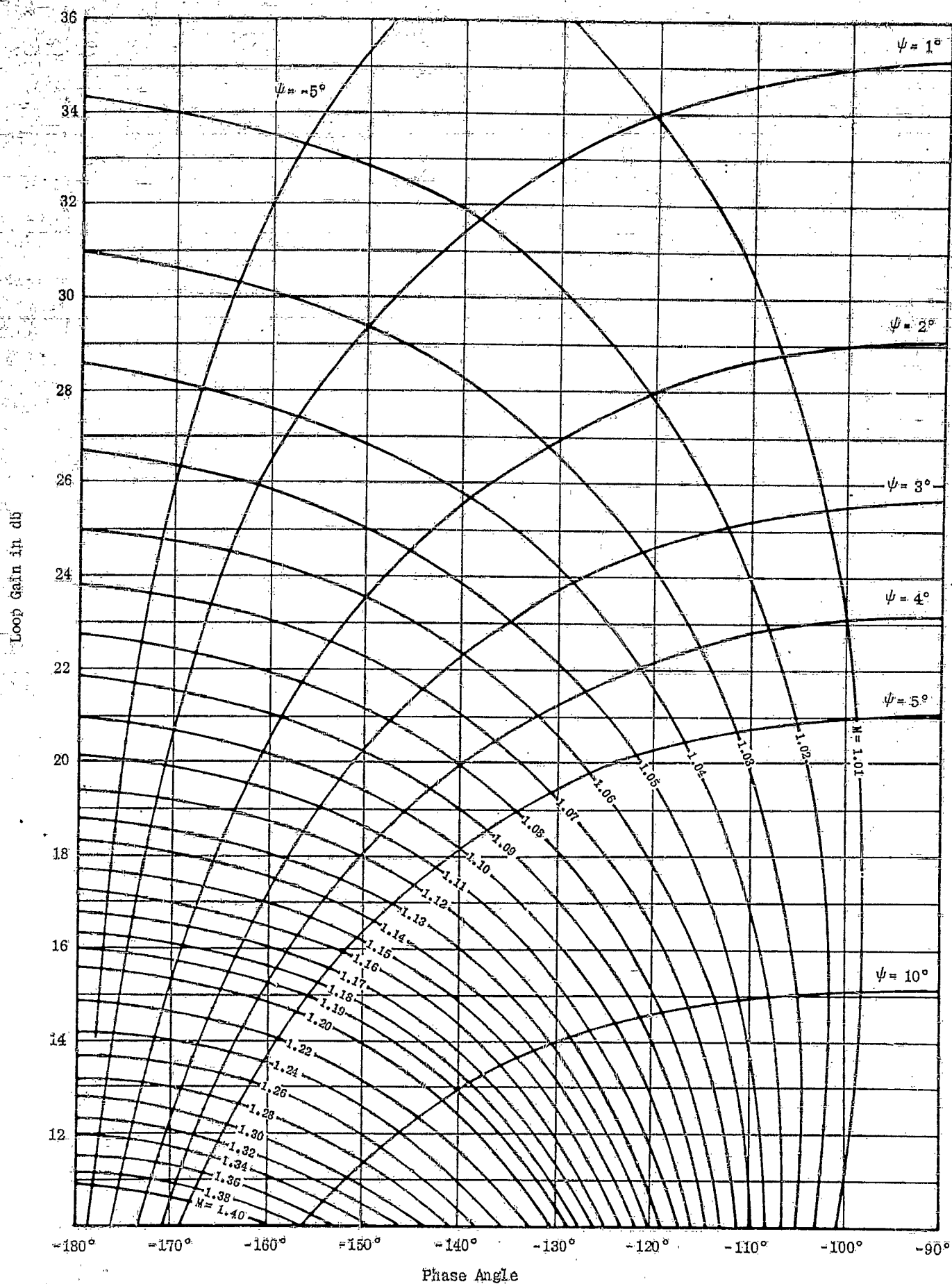
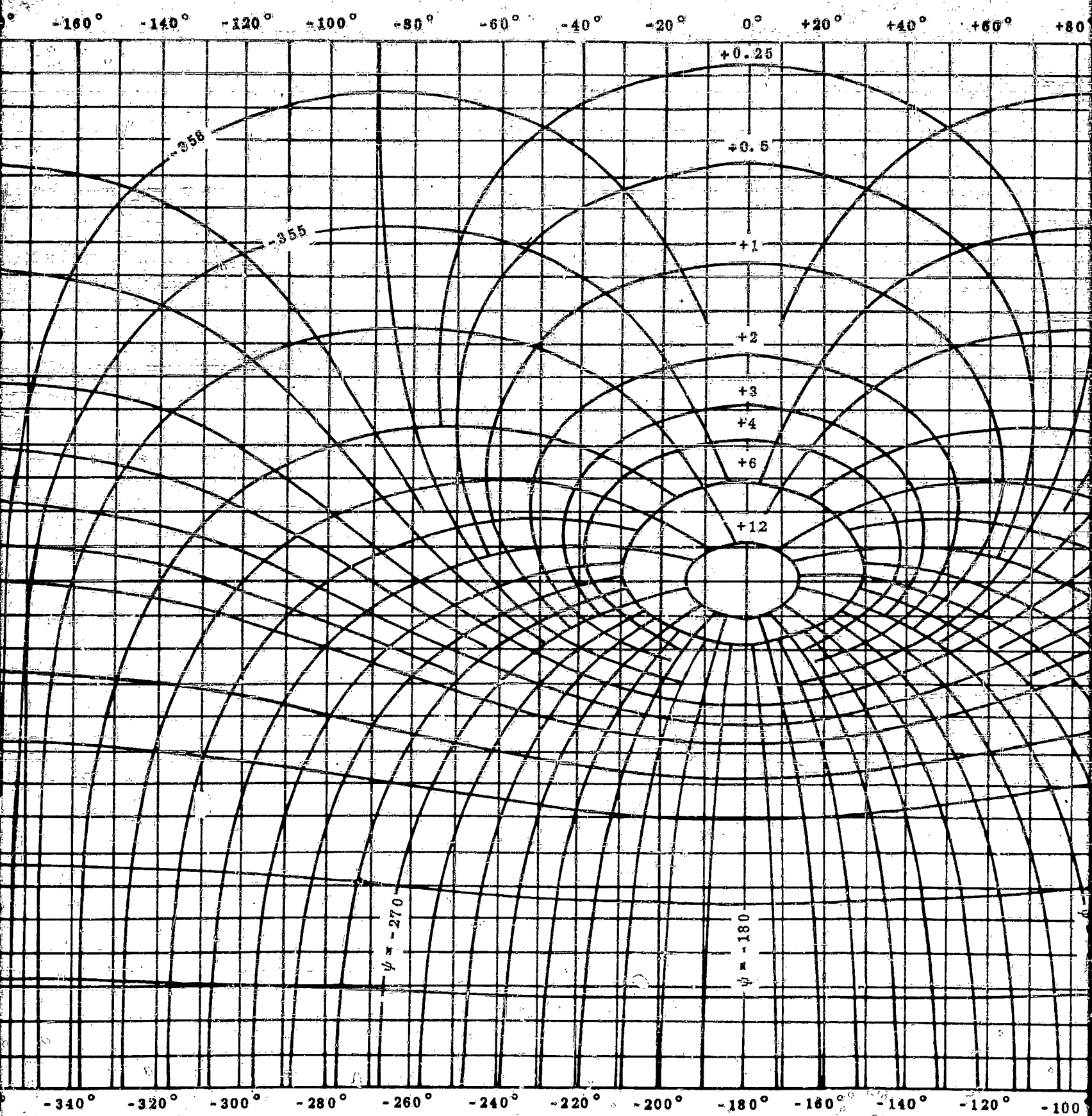
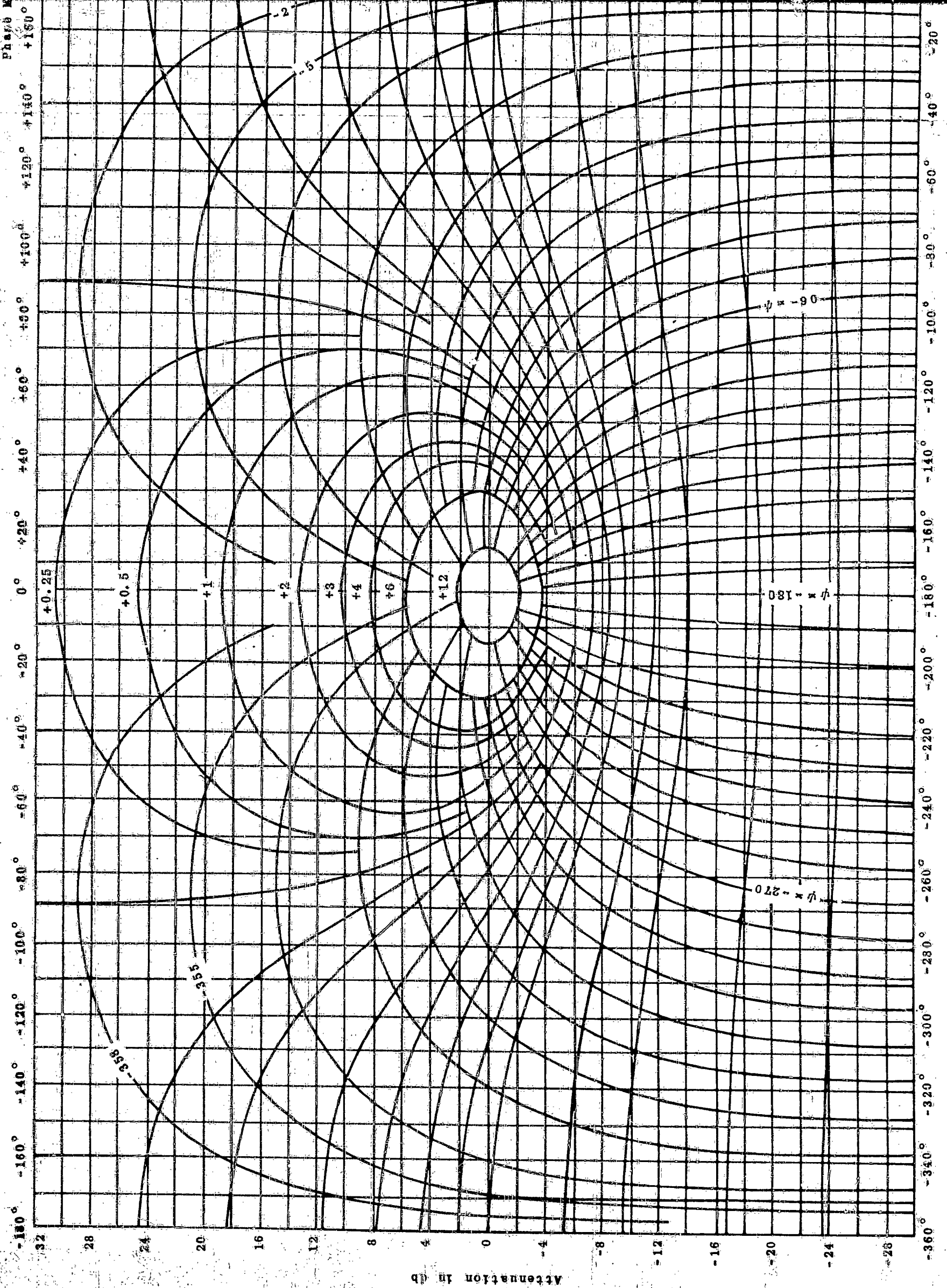


Figure A-18. Nichols Chart



Phase



Attenuation in db

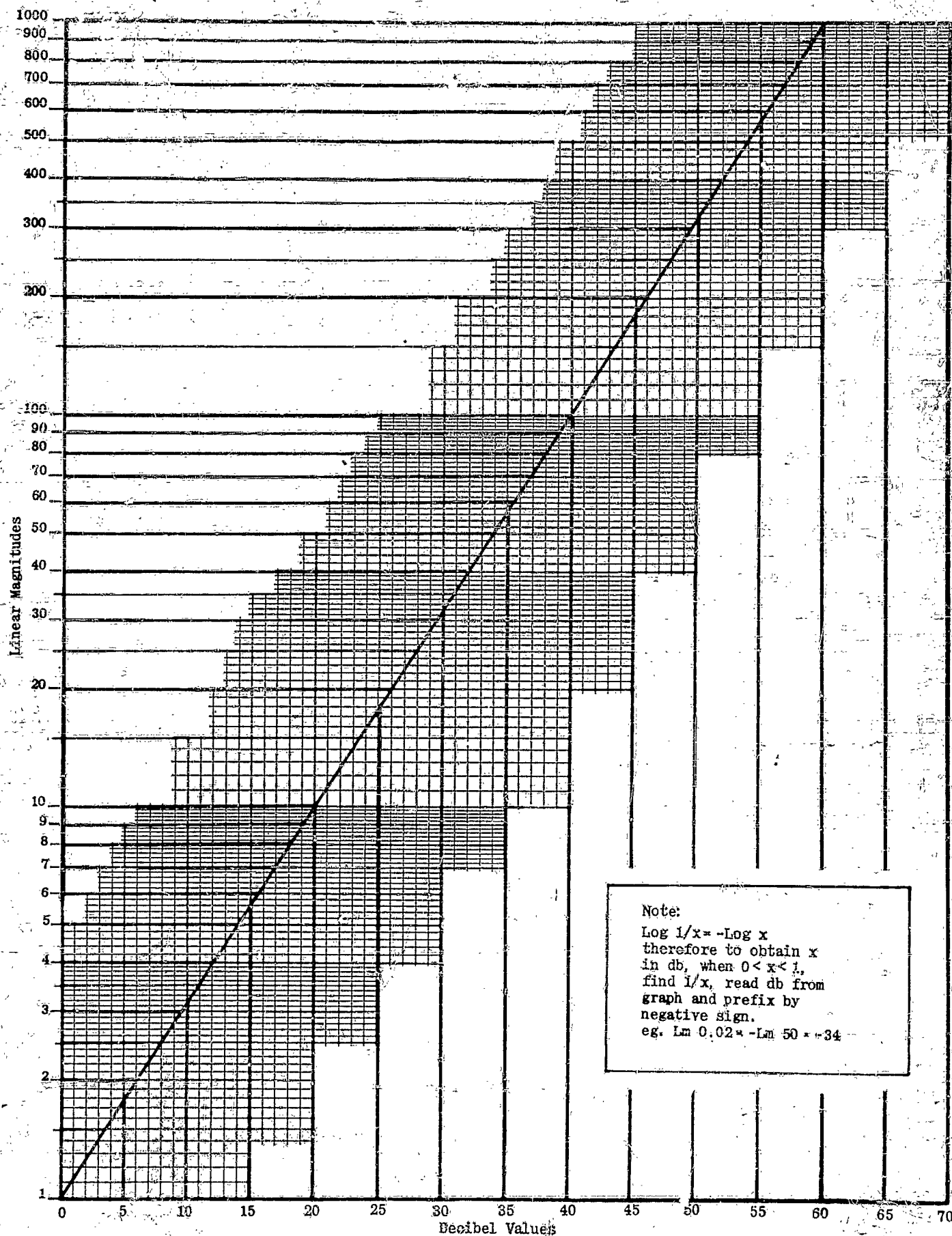


Figure A-20. Linear Magnitudes vs. Decibel Values

SECTION A I - ROOTS OF ALGEBRAIC EQUATIONS

One method of determining the transient response of a dynamical system is to calculate the roots of the characteristic equation (the characteristic equation is the denominator polynomial of the transfer function set equal to zero). As was pointed out in the main body of the text, most of the methods available to find roots are tedious and are avoided when possible. However some simple cases may be effectively handled in this way, and for this reason the following discussion is presented.

The problem is: Given the equation,

$$(A I-1) \quad z^n + a_1 z^{n-1} + \dots + a_{n-1} z + a_n = 0$$

find the values of z for which the equation is satisfied, when the $a_i (i = 1, 2, 3, \dots, n)$ are known real quantities.

It can be proved rigorously that the desired roots are never anything other than complex numbers, with real numbers considered special cases of complex numbers.

It is a matter of common experience that real roots are much more easily located than are the complex ones; it is also true that some of the methods for finding all of the roots tend to work out considerably better when the equation has only complex roots. A logical procedure is then:

1. Determine real roots and remove them from characteristic equation.
2. Determine complex roots from remainder.

Standard graphical and numerical methods for doing this are well known, such as the Newton-Raphson method, Horner's method, and the so-called "method of false position." For those wishing to study these in detail, standard references* provide ample material. They all depend, however, on first obtaining at least a rough approximation to a real root, and then improving the approximation. This can be done to any degree of accuracy desired, by these methods.

In the general case, there is no information about the characteristic equation to indicate if there are any real roots. The initial step is to obtain some information on this point. A method due to J.C. F. Sturm, a French mathematician of the early nineteenth century, which determines if an equation has real roots, will now be presented. The method may be stated rather simply and is presented here without proof.

Let the characteristic equation be

$$P(z) = a_0 z^n + a_1 z^{n-1} + \dots + a_n = 0$$

From the given $P(z)$, a sequence of functions is derived as follows:

The first Sturm function, referred to as that of zero

order, is $P_0(z)$ itself:

$$(A I-2) \quad S_0 = P(z)$$

The second Sturm function is the derivative of P with respect to z , dP/dz .

$$(A I-3) \quad S_1 = \frac{dP}{dz}$$

The third Sturm function is derived from S_0 and S_1 , as follows: Divide S_0 by S_1 ; the result will be a quotient Q_1 and a remainder R_1 ; the latter will be of degree $n-2$. The quotient is of no further use in the process, and is discarded. The remainder, with changed algebraic sign, is S_2 . In symbols, if $S_0/S_1 = Q_1 + (R_1/S_1)$ then:

$$(A I-4) \quad S_2 = -R_1$$

The further Sturm functions of the sequence are obtained by using the following formula: $\frac{S_{k-1}}{S_k} = Q_k + \frac{R_k}{S_k}$; and the definition

$$(A I-5) \quad S_{k+1} = -R_k$$

that is, by repetitions of the general process giving R_1 .

Since the degree of each successive Sturm function is one less than that of the preceding one, the process will finally end with an S_n which is a constant, not dependent upon z .

The whole sequence of Sturm functions for $P(z)$ is then:

$$(A I-6)$$

$$\begin{aligned} S_0 &= z^n + a_1 z^{n-1} + \dots + a_{n-1} z + a_n \\ S_1 &= n z^{n-1} + (n-1) a_1 z^{n-2} + \dots + 2 a_{n-2} z + a_{n-1} \\ S_2 &= -\frac{1}{n} \left[2 a_2 - \frac{(n-1) a_1^2}{n} \right] z^{n-2} - \frac{1}{n} \left[3 a_3 - \frac{(n-2) a_1 a_2}{n} \right] z^{n-3} \\ &\quad - \dots - \frac{1}{n} \left[(n-1) a_{n-1} - \frac{2 a_1 a_{n-2}}{n} \right] z - \frac{1}{n} \left[n a_n - \frac{a_1 a_{n-1}}{n} \right] \\ &\quad \vdots \\ S_n &= K \end{aligned}$$

where K , of course, is a constant.

The location of the real roots of $P(z)$ is then carried out by applying the following rule:

If it is desired to determine the number of real roots of $P(z) = 0$ lying between the (real) values a and b of z , evaluate each of the functions $S_0(a), S_1(a), \dots, S_n(a)$ and $S_0(b), S_1(b), \dots, S_n(b)$. The number of real roots between a and b is then the difference in the number of changes of sign between the two sequences of functions.

It is to be noted that since only the algebraic signs of the functions of the two sequences is of interest in this rule, the arithmetic used in establishing what the sequences are can be considerably simplified. Any function of the sequence (A I-6), for example, can be divided through by any positive constant which will

* Willers, Dr. F.A., *Practical Analysis*, Dover, N.Y. 1948
Milne, W.E., *Numerical Calculus*, Princeton U. Press, Princeton, N.J., 1948
Scarborough, J.B., *Numerical Mathematical Analysis*, The John Hopkins Press, Baltimore, Md. 1930

make the work easier. Instead of using S_1 as it stands in (A.1-6), it may be divided through by n to facilitate computation.

The rule makes it very easy to settle the question of whether $P(z)$ has any real roots or not. In this case, the constants a and b may be taken as $+\infty$ and $-\infty$, respectively, and the evaluation of the functions of the sequences is then simply a matter of inspection of the algebraic sign of the terms of highest degree in z . If the S_k have two changes of sign for $z = +\infty$ and four changes of sign for $z = -\infty$, then, by the rule, there are exactly two real roots of the characteristic equation $P(z) = 0$.

It should also be mentioned that the presence of positive real roots may be determined very easily by this method. In this case, the functions are to be evaluated for $z = +\infty$ and $z = 0$. The former gives simply the algebraic sign of the terms of highest degree, as before; and the latter, that of the constant terms; examination of this range of the variable can then also be carried out by inspection.

The use of the Sturm functions can be extended to the process of getting initial approximate values of the roots for refinement by one of the standard methods. The a and b used above in the rule are any real numbers at all, and by taking trial values of these, the Sturm process can be used to locate ranges of any size at all in which roots must lie. Theoretically, the Sturm functions can be used to locate the roots with any accuracy desired; however, the computational work involved in their use for this purpose is considerably greater than for more usual techniques, and their greatest utility is achieved in merely determining fairly broad ranges of values in which the roots must lie.

An important exception to the use of these functions must now be noted; the Sturm functions will always do the job they are supposed to, when they exist; however, if S_0 contains higher order roots, it will be found that the result of the first division, S_0/S_1 , has no remainder. In this case, the process does not apply. There are ways of avoiding this, but in controls work the case of multiple roots of $P(z)$ occurs so seldom that they are not worth considering here.

For the sake of illustration, a numerical application of Sturm functions now follows; consider the characteristic equation:

$$(A.1-7) \quad z^6 - z^5 - 14z^4 - z^3 + 25z^2 + 38z + 24 = 0$$

The corresponding sequence of Sturm functions may be computed as:

$$(A.1-8) \quad S_0 = z^6 - z^5 - 14z^4 - z^3 + 25z^2 + 38z + 24$$

$$S_1 = z^5 - 0.8332z^4 - 9.3332z^3 - 0.5z^2 + 8.3332z + 6.333$$

$$S_2 = z^4 + 0.4282z^3 - 3.4502z^2 - 6.8792z - 5.214$$

$$S_3 = z^3 - 0.3802z^2 - 0.910z - 0.045$$

$$S_4 = z^2 + 2.772z + 2.351$$

$$S_5 = -z - 1.361$$

$$S_6 = -0.431$$

The information obtainable from these, in accordance with the rule given above, is best indicated in tabular form:

	S_0	S_1	S_2	S_3	S_4	S_5	S_6	Number of sign changes
$+\infty$	+	+	+	+	+	-	-	1
0	+	+	-	-	-	-	-	3
$-\infty$	+	-	-	+	+	+	-	5

The crosses have been placed between the signs of the functions of the sequences to indicate where the sign changes occur.

It is seen that there is a difference of four changes of sign between $-\infty$ and $+\infty$; the conclusion is that the given $P(z)$ has exactly four real roots; since there are differences in changes of sign of two between $-\infty$ and 0, and two between 0 and $+\infty$, there are two positive and two negative real roots. Actually, this equation was arrived at by multiplying out the expression: $(S+1)(S-2)(S+3)(S-4)(S^2+S+1)$; it may be seen that this has exactly the distribution of roots predicted by the Sturm function process.

Once the total number of real roots, and their approximate locations are known, any of the standard methods may be applied to obtain better approximations to the roots. When these have been carried out to whatever accuracy is desired, the degree of the equation may then be successively decreased by one by dividing out each of these roots.

The process of synthetic division is particularly recommended in this determination of the real roots; if r is a root of the equation, synthetic division by it leaves the remainder zero, and the partial quotients are the coefficients of the equation of next lower degree. If r is not an exact root (to that order of approximation desired) the last figure resulting from the division is the remainder, and may be used as a measure of how far away from the true root is the approximation which has just been tried. A closer approximation can then be obtained by any of the standard methods, and the division repeated with the new trial divisor.

When all the real roots have been removed, not all the factors of the original equation may be accounted for. It is then necessary to have a process which will extract complex roots. Several methods for this purpose are in more or less frequent use; of these, the most common and best seem to be Graeffe's root squaring method, and Lin's method.

Graeffe's method suffers from certain defects not found in some of the others. In the first place, the numerical work involved is not of an iterative nature, so that an error in computation makes all the rest of the application of the process worthless. It is thus essential to check all roots obtained by substitution into the original equation. Also, the roots are not obtained initially as their correct values, but as the second, fourth, or some higher even power of the roots of the original equation. This means that it is necessary to extract the square, fourth, ..., (2ⁿ)th root of the complex number which the Graeffe method gives. This gives rise to additional possibilities

of error in the computations. For these reasons, as well as the fact that the process is rather commonly known, this method will not be discussed further here.

Lin's method is probably the most usable of all the present processes for the extraction of complex roots. It takes out the roots as quadratic factors; since each pair of complex roots gives rise to one such factor, it may be said to yield the roots themselves directly, in a form requiring only the application of the quadratic formula. It has the further advantage that the computation is iterative; errors tend to slow down the work, but do not absolutely prevent getting the correct final result.

It suffers from what appears to be the common defect of all present methods of accomplishing the same end; if two pairs of complex roots are very nearly equal in magnitude, the process converges slowly, and may need many repetitions of its basic procedure before a sufficiently accurate result can be obtained. It is also true that in some rare cases Lin's method diverges rather than converges. The conditions under which this may occur are not fully understood; it happens only very seldom, in any event, and is practically always obviated by taking out all real roots first.

Essentially, Lin's method is this: The last three terms of a polynomial are used as a first approximation to a quadratic factor. The polynomial is divided by the first approximation. The result is a quotient and a remainder.

Certain coefficients in the quotient are used to get a new trial divisor from the old one.

In this process, the remainder has no use but to serve as a measure of the closeness of the approximation to the correct value. The method is repeated until the remainder is within some assigned limit of approximation; or, until it remains of the same magnitude after several repetitions of the process. If this last occurs, it indicates that the result is as accurate as can be attained with the number of digits used in the coefficients.

It should be noted here that the first cycle of Lin's process gives an approximate factorization of the polynomial. If the coefficients are literal, this first cycle is all that it is practicable to use. However, in practical situations as in the case of the stability quartics, the approximate factorization works out excellently.

Figure A I-1 shows the application of the first cycle of Lin's method to a general case. This illustration shows how the first trial divisor is obtained from the quadratic terms of the polynomial, in line (2); it also illustrates the notation used in the quotient, and shows in line (4) how the next trial divisor is obtained from the appropriate terms of the quotient and dividend. The subsequent trial divisors are obtained in a similar manner.

If the polynomial being investigated has a pair of complex roots of very nearly equal magnitude, the con-

$$(1) \quad z^n + a_{n-1}z^{n-1} + \dots + a_2z^2 + a_1z + a_0 = 0$$

(2) First trial divisor

$$z^2 + \frac{a_1}{a_2} z + \frac{a_0}{a_2} \quad d_1^{(1)} \triangleq \frac{a_1}{a_2}; \quad d_0^{(1)} \triangleq \frac{a_0}{a_2}$$

$$(3) \quad z^2 + d_1^{(1)}z + d_0^{(1)} \quad z^n + a_{n-1}z^{n-1} + \dots + a_2z^2 + a_1z + a_0 \quad (z^{n-2} + (a_{n-1} - d_1^{(1)})z^{n-3} \\ + \dots + d_2^{(1)}z^2 + q_1^{(1)}z + q_0^{(1)}) \\ \hline (a_{n-1} - d_1^{(1)})z^{n-1} + \dots \\ \hline (a_{n-1} - d_1^{(1)})z^{n-1} + \dots$$

$$\frac{a_1 z + a_0}{(r_1^{(1)} z + r_0^{(1)})} \quad \text{remainder}$$

(A) Second trial division

$$d_1^{(2)} = \frac{a_1 q_0^{(1)} - a_0 q_1^{(1)}}{[q_1^{(1)}]^2}$$

$$d_0^{(2)} = \frac{a_0}{d_0^{(1)}}$$

Figure AI-1. Lin's Method for the Complex Roots of a Polynomial

Appendix
Section A

$$8 + 7 + 5 = 8 (1 + 7/8 + 5/8) = 8 (1 + 0.875 + 0.625)$$

$$1 + 0.875 + 0.625) \quad 1 + 3.000 + 8.000 + 7.000 + 5.000 \quad (1 + 2.125 + 5.516)$$

$$-1 - 0.875 - 0.625$$

$$2.125 + 7.375 + 7.000$$

$$-2.125 - 1.859 - 1.328$$

$$5.516 + 5.672 + 5.000$$

$$-5.516 - 4.627 - 3.448$$

$$+ 1.845 + 1.552$$

$$a_1 = 7$$

$$q_1 = 2.125$$

$$a_0 = 5$$

$$q_0 = 5.516$$

$$d_1 = \frac{a_1 q_0 - a_0 q_1}{q_0^2} = \frac{27.987}{30.426256} = 0.920$$

$$d_0 = \frac{a_0}{q_0} = \frac{5}{5.516} = 0.906$$

$$1 + 0.920 + 0.906) \quad 1 + 3.000 + 8.000 + 7.000 + 5.000 \quad (1 + 2.080 + 5.180)$$

$$-1 - 0.920 - 0.906$$

$$2.080 + 7.094 + 7.000$$

$$-2.080 - 1.914 - 1.884$$

$$5.180 + 5.116 + 5.000$$

$$-5.180 - 4.766 - 4.693$$

$$+ 0.350 + 0.307$$

$$d_1 = \frac{(7)(5.180) - (5)(2.080)}{(5.180)^2}$$

$$= \frac{25.86}{26.8324} = 0.964$$

$$d_0 = \frac{5}{5.18} = 0.965$$

$$1 + 0.964 + 0.965) \quad 1 + 3.000 + 8.000 + 7.000 + 5.000 \quad (1 + 2.036 + 5.072)$$

$$-1 - 0.964 - 0.965$$

$$2.036 + 7.035 + 7.000$$

$$-2.036 - 1.963 - 1.965$$

$$+ 5.072 + 5.035 + 5.000$$

$$-5.072 - 4.889 - 4.894$$

$$0.146 + 0.106$$

$$d_1 = \frac{(7)(5.072) - (5)(2.036)}{(5.072)^2}$$

$$= \frac{25.324}{25.725184} = 0.984$$

$$d_0 = \frac{5}{5.072} = 0.986$$

$$1 + 0.984 + 0.986) \quad 1 + 3.000 + 8.000 \quad (1 + 2.016 + 5.030)$$

$$-1 - 0.984 - 0.986$$

$$2.016 + 7.014 + 7.000$$

$$-1.984 - 1.988$$

$$5.030 + 5.012 + 5.000$$

$$-4.950 - 4.960$$

$$0.062 + 0.040$$

$$a_1 = \frac{(7)(5.030) - (5)(2.016)}{(5.030)^2}$$

$$= \frac{25.130}{25.3009}$$

$$= 0.993$$

$$d_0 = \frac{5}{5.03} = 0.994$$

$$1 + 0.993 + 0.994) \quad 1 + 3.000 + 8.000 \quad (1 + 2.007 + 5.013)$$

$$2.007 + 7.006 + 7.000$$

$$-1.993 - 1.995$$

$$5.013 + 5.005 + 5.000$$

$$-4.978 - 4.983$$

$$0.027 + 0.017$$

Figure A1-2. Example of Lin's Method

vergence depends directly upon the relative magnitudes of the roots. For this reason, when the process appears to be going too slowly, it is worth while to subject the polynomial to a transformation which will increase differences in magnitudes of the roots. Graeffe's method may be used for this. If there is any appreciable separation in the size of the roots, squaring them once or twice will give a transformed equation for which Lin's process will converge with satisfactory rapidity. As a matter of fact, the roots are usually far enough apart in the equations met with in controls work so that it is seldom necessary to transform the polynomial in this way.

As a numerical example, the location of the roots of $(z^2 + 2z + 5)(z^2 + z + 1) = z^4 + 3z^3 + 8z^2 + 7z + 5 = 0$ will be obtained from a knowledge of the untransformed form.

SECTION A II - THE ROUTH-HURWITZ STABILITY CRITERION

In dealing with stability problems one oftentimes has to decide if the roots of a given algebraic equation of n^{th} degree, such as

$$(A II-1) \quad a_0 z^n + a_1 z^{n-1} + \dots + a_{n-1} z + a_n = 0$$

all have negative real parts. The coefficients $a_0 \dots a_n$ are real numbers.

A general criterion given by Hurwitz is that the n determinants:

$$(A II-2) \quad D_1 = a_1; \quad D_2 = \begin{vmatrix} a_1 & a_0 \\ a_3 & a_2 \end{vmatrix}; \quad D_3 = \begin{vmatrix} a_1 & a_0 & 0 \\ a_3 & a_2 & a_1 \\ a_5 & a_4 & a_3 \end{vmatrix}$$

$$\dots D_n = \begin{vmatrix} a_1 & a_0 & 0 & \dots & 0 \\ a_3 & a_2 & a_1 & \dots & 0 \\ a_5 & a_4 & a_3 & \dots & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot \\ a_{2n-1} & a_{2n-2} & a_{2n-3} & \dots & a_n \end{vmatrix}$$

must all be positive. The coefficient a_0 must be positive or must be made positive (by multiplying the equation by -1).

This stability criterion results in simple rules for quadratic, cubic and fourth order equations

For the above three types of equations the rule states that all coefficients must be positive. In addition, for the cubic equation

$$(A II-3) \quad a_0 z^3 + a_1 z^2 + a_2 z + a_3 = 0$$

there is the condition that

$$(A II-4) \quad a_1 a_2 - a_0 a_3 > 0$$

For the equation of fourth degree

$$(A II-5) \quad a_0 z^4 + a_1 z^3 + a_2 z^2 + a_3 z + a_4 = 0$$

Figure A II-2 shows the necessary work.

It will be noted that four cycles of Lin's process determines the quadratic factors with errors of less than 1% in their coefficients.

The magnitudes of the roots of the factors are 1 and $\sqrt{5} = 2.24$. It will be seen that the ratio of the magnitudes is not very much greater than unity; despite this, the convergence of the process is fairly rapid.

Detached coefficients have been used in the working out of this example to shorten the process. A little practice in applying the method will enable one to simplify it still further by omitting the writing of various numbers occurring in the process. However, it was felt that the whole work should be shown here to facilitate checking one's understanding of the process.

the additional condition is that

$$(A II-6) \quad a_2(a_1 a_2 - a_0 a_3) - a_1^2 a_4 > 0$$

For equations of higher degree than the fourth all the determinants (A II-2) must be formed and the value of each must be tested for its sign.

Another, more general criterion, the Routh criterion, is given in chapter VII, section 9, of "Transients in Linear Systems" by Gardner and Barnes. In addition to indicating instability it also discloses the number of roots with positive real parts as well as the number of purely imaginary roots.

As an example of application of the above stability criterion let it be required to investigate the stability of the motion of the pendulum shown below (see figure A II-1). The pendulum is rotated about its axis $o-o$ at an angular velocity ω rad./sec.; it is deflected laterally thru a small angle α ; will this angular deflection die out, or will it diverge?

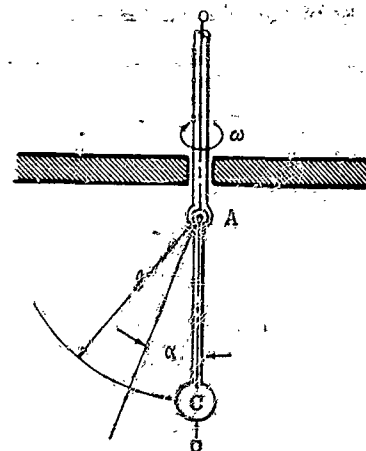


Figure A II-1. Pendulum (Stability)

It is assumed that ω is constant; that the mass of the pendulum is concentrated in the bob, c ; and that the lateral motion is opposed by viscous friction.

Appendix Section A III

The equation of motion of this system is:

$$(A II-7) \quad \ddot{x} + 2n\dot{x} + \left(\frac{g}{l} - \omega^2\right)x = 0$$

where $2n$ is the damping coefficient; g is the gravitational acceleration constant. Assuming a solution of (A II-7) of the form $x = e^{st}$ there results the equation

$$(A II-8) \quad s^2 + 2ns + \left(\frac{g}{l} - \omega^2\right)s = 0$$

and

$$(A II-9) \quad s_{1,2} = -n \pm \sqrt{n^2 - \left(\frac{g}{l} - \omega^2\right)}$$

If $g/l - \omega^2 > 0$, the real parts of both roots, s_1 and s_2 , are negative; this corresponds to a condition of stability, i.e. the oscillation will die out; if $g/l - \omega^2 < 0$, or, calling $g/l - \omega^2 = -p$,

$$(A II-10) \quad s^2 + 2ns - p = 0$$

it is seen that the root $s = -n + \sqrt{n^2 + p}$ is positive; this corresponds to instability of the lateral motion.

This shows that (A II-10), which does not comply with the condition that all coefficients be positive, indicates an unstable motion.

SECTION A III -- FUNCTION OF A COMPLEX VARIABLE

(A pre-requisite for the understanding of this section is a working knowledge of the elements of the theory of complex numbers. This includes the properties of complex numbers, the various ways of representing them algebraically and graphically; and the fundamental operations of addition, subtraction, multiplication, division, evolution and involution of complex numbers.)

CONTENTS OF THIS SECTION

This section introduces the concept of functions of complex variables; it defines analytic functions and discusses some of their properties; it introduces the complex w -plane, in which functions w of the complex variable z are plotted; it discusses poles and zeros of w ; finally it deals with the line integral of w .

The material presented here is a useful tool in the study of various technical subjects, especially in the study of electrical circuits and of certain phases of Servomechanism Theory.

INTRODUCTION OF $w(z)$, A FUNCTION OF THE VARIABLE z

Let z be a complex number $z = x + iy$. When $y = 0$, z becomes real. Thus real numbers may be considered a special case of complex numbers.

The function $y = f(x)$ of the real variable x may be represented as a curve in the xy -plane. The value of y which corresponds to a given x is found on the same horizontal line as the point on the curve which

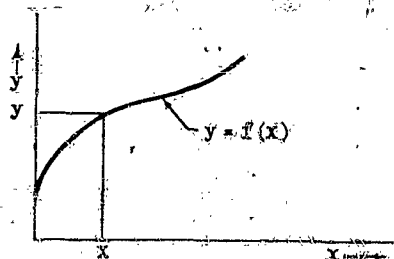


Figure A III-1. Real Function

lies vertically above the given x . It may be stated that: As x moves along the x -axis, y moves along the curve $f(x)$ as shown in figure A III-1.

Now, one may think of the function $y = f(x)$ of the real variable x as a special case of the function $w = F(z)$ of the complex variable z . The function $w = F(z)$ is a generalization of the function $y = f(x)$, just as z is a generalization of x .

Since $z = x + iy$, $w = F(z) = F(x + iy)$. It is of interest to find the relationship between w , and x and y , when the relationship between w and z is known.

The following example shows how this can be done. Let

$$(A III-1) \quad w = f(z) = z^2$$

Then

$$(A III-2) \quad w = (x + iy)^2 = (x^2 - y^2) + i 2xy$$

Thus, w is itself a complex number; designating the real part of w by u and the imaginary part by v ,

$$(A III-3) \quad u = x^2 - y^2 \quad v = 2xy$$

and

$$(A III-4) \quad w = u + iv$$

It may be seen that the real and also the imaginary part are functions of x and y .

$$(A III-5) \quad w = F(x, y) + i\phi(x, y)$$

Another example; let

$$(A III-6) \quad w = \frac{1}{z}$$

Then

$$(A III-7) \quad w = \frac{1}{x + iy} = \frac{x - iy}{x^2 + y^2} = \frac{x}{x^2 + y^2} - i \frac{y}{x^2 + y^2} = u + iv$$

where

$$(A III-8) \quad u = \frac{x}{x^2 + y^2}; \quad v = \frac{-y}{x^2 + y^2}$$

* z is used as the symbol for a complex number (instead of s as in the body of this volume) to correspond with traditional mathematical usage.

THE DERIVATIVE OF w WITH RESPECT TO z .

Since w is a function of z , it is of interest to find the derivative dw/dz .

The real function $f(x) = f(x)$ has a derivative dy/dx (see figure A III-2); this is a scalar quantity for any given value of x ; it represents the magnitude of the slope of the curve $f(x)$ at the given point x . The expression for dy/dx may contain, in addition to x and functions of x , also functions of y ; but since y must lie on the known curve $y = f(x)$, y is completely specified for a given x .

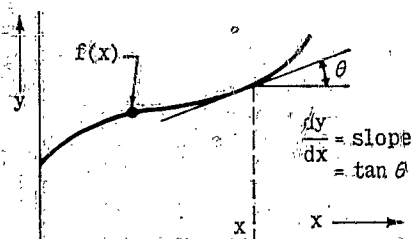


Figure A III-2. Derivative of a Real Function

The derivative dw/dz will now be found.

The derivative dy/dx has been defined as the limit of the ratio $\Delta y/\Delta x$ as the increment Δx approaches zero. This Δx is a change in x along the x -axis and Δy is the corresponding change in y as the latter moves along the curve (see figure A III-3).

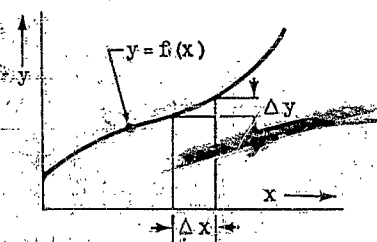


Figure A III-3. Finite Change in y

If one were to define dw/dz as the limit of the ratio $\Delta w/\Delta z$ as the vector z is given an increment Δz , and Δz approaches zero, a difficulty would arise; for while, in the case of $\Delta y/\Delta x$, x is constrained to move thru the increment Δx along the x -axis, z can "move through Δz " in any direction at all; and this would mean that dw/dz might have any number of values, depending on the direction of Δz . (See figure A III-4).

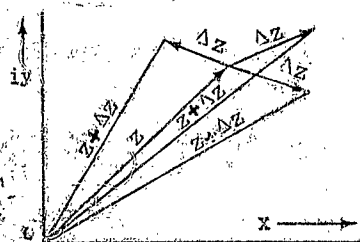


Figure A III-4. Possible Changes of z

Obviously, in order to have a definite answer for a derivative and still have a definite single answer, the group of functions to be dealt with must be narrowed down to those which are so constituted that dw/dz has the same value regardless of the direction of Δz . Just what the condition is that is imposed by this requirement can be found by the following reasoning: Let

$$(A \text{ III-9}) \quad \frac{dw}{dz} = \lim_{\Delta z \rightarrow 0} \frac{\Delta w}{\Delta z} = \lim_{\Delta z \rightarrow 0} \frac{\Delta u + i \Delta v}{\Delta x + i \Delta y}$$

and let Δz lie along the x -axis; i.e., $\Delta z = \Delta x$; then

$$(A \text{ III-10}) \quad \frac{dw}{dz} = \lim_{\Delta x \rightarrow 0} \frac{\Delta u + i \Delta v}{\Delta x} = \lim_{\Delta x \rightarrow 0} \left(\frac{\Delta u}{\Delta x} + i \frac{\Delta v}{\Delta x} \right) = \frac{\partial u}{\partial x} + i \frac{\partial v}{\partial x}$$

Now, let Δz lie along the y -axis; i.e., $\Delta z = i \Delta y$; then

$$(A \text{ III-11}) \quad \frac{dw}{dz} = \lim_{\Delta y \rightarrow 0} \frac{\Delta u + i \Delta v}{i \Delta y} = \lim_{\Delta y \rightarrow 0} \left(\frac{\Delta u}{i \Delta y} + \frac{\Delta v}{\Delta y} \right) = -i \frac{\partial u}{\partial y} + \frac{\partial v}{\partial y}$$

If the derivative is to be independent of the direction of Δz , these two expressions for dw/dz must be identically the same; hence,

$$(A \text{ III-12}) \quad \frac{\partial u}{\partial x} + i \frac{\partial v}{\partial x} = -i \frac{\partial u}{\partial y} + \frac{\partial v}{\partial y}$$

or;

$$(A \text{ III-13}) \quad \frac{\partial u}{\partial x} = \frac{\partial v}{\partial y} \quad \text{and} \quad \frac{\partial v}{\partial x} = -\frac{\partial u}{\partial y}$$

A function $w = f(z) = f(u, v) = F(x, y) + i\phi(x, y)$ which satisfies the condition (A III-13) is called "analytic" at the point z where it has a derivative. The equations (A III-13) are known as the Cauchy-Riemann conditions.

For the function $w = z^2$ given above:

$$(A \text{ III-14}) \quad \begin{aligned} u &= x^2 - y^2 & \left\{ \begin{aligned} \frac{\partial u}{\partial x} &= 2x \\ \frac{\partial u}{\partial y} &= -2y \end{aligned} \right. \\ v &= 2xy & \left\{ \begin{aligned} \frac{\partial v}{\partial x} &= 2y \\ \frac{\partial v}{\partial y} &= 2x \end{aligned} \right. \end{aligned}$$

This function is analytic at any point in the finite z plane. It satisfies the condition (A III-13).

Thus, $dw/dz = d(z^2)/dz = 2z$; the differentiation is carried out as in the case of a real function; the derivative, $2z$, has a definite value of any point z .

For the function $w = 1/z$:

$$(A \text{ III-15}) \quad \begin{aligned} u &= \frac{x}{x^2 + y^2} & v &= -\frac{y}{x^2 + y^2} \\ \frac{\partial u}{\partial x} &= \frac{(x^2 + y^2) - x(2x)}{(x^2 + y^2)^2} = \frac{y^2 - x^2}{(x^2 + y^2)^2} \\ \frac{\partial u}{\partial y} &= \frac{-(x^2 + y^2)(2y) - y^2(2y)}{(x^2 + y^2)^3} = \frac{-2xy^2 - 2y^3}{(x^2 + y^2)^3} = \frac{-2y(x^2 + y^2)}{(x^2 + y^2)^3} = \frac{-2y}{(x^2 + y^2)^2} \\ \frac{\partial v}{\partial x} &= \frac{0 - y(2x)}{(x^2 + y^2)^2} = \frac{-2xy}{(x^2 + y^2)^2} \\ \frac{\partial v}{\partial y} &= \frac{-(x^2 + y^2) - y(2y)}{(x^2 + y^2)^2} = \frac{-x^2 - y^2 - 2y^2}{(x^2 + y^2)^2} = \frac{-x^2 - 3y^2}{(x^2 + y^2)^2} \end{aligned}$$

Appendix Section A III

Thus, $w = 1/z$ is analytic at all points z , except $z = 0$, at which w and dw/dz do not exist, nor does the derivative. At all other points

$$(A III-16) \quad \frac{dw}{dz} = \frac{d(1/z)}{dz} = -\frac{1}{z^2}$$

It can be shown that $w = z^n$ and $dw/dz = nz^{n-1}$ exist and (n-integral) are analytic at all points in the z -plane when $n > 0$. When $n < 0$, w has a derivative at all points except at the origin, i.e., at $z = 0$.

The sum or the product of two analytic functions is also an analytic function. Thus, any polynomial

$$(A III-17) \quad w = a_0 + a_1 z + a_2 z^2 + \dots + a_n z^n$$

is an analytic function, in the finite part of the z -plane.

The quotient of two such polynomials is analytic except at those values of z which make the denominator vanish, i.e., for which w is infinite.

It can also be shown that $\sin z$, $\cos z$, and other trigonometric functions of z , all of which are expressible in power series of z , are analytic except at certain points; and the same is true of the hyperbolic functions of z . Also: if $\sin z$ is an analytic function of z , and $(1+z^2)^2$ is also analytic, then $\sin(1+z^2)$ is also analytic; i.e., it can be shown that an analytic function of an analytic function is itself analytic. And in all cases it will be found that the analytic functions satisfy the Cauchy-Riemann conditions (A III-13).

An example of a non-analytic function is

$$(A III-18) \quad w = z \bar{z} = (x + iy)(x - iy) = x^2 + y^2 = u + iv$$

where $u = x^2 + y^2$, $v = 0$; obviously $\partial u / \partial x \neq \partial v / \partial y$.

It should also be mentioned that a function fails to be analytic at a point where it is multiple-valued. An example is $w = z^{1/2} = \pm \sqrt{z}$; for this simple function, things can be redefined, of course, so that the plus of the double sign is always taken, which makes $z^{1/2}$ single valued; but until this is done the function as it stands is not analytic.

From this point on only analytic functions will be dealt with.

EXPONENTIAL FUNCTIONS OF z

What is the meaning of $w = e^z = e^{x+iy}$? To arrive at the answer it is convenient to consider the series expansion of $\cos \theta$, $\sin \theta$ and e^θ ;

$$(A III-19) \quad \cos \theta = 1 - \frac{\theta^2}{2!} + \frac{\theta^4}{4!} - \frac{\theta^6}{6!} + \dots$$

$$\sin \theta = \theta - \frac{\theta^3}{3!} + \frac{\theta^5}{5!} - \frac{\theta^7}{7!} + \dots$$

$$e^\theta = 1 + \frac{\theta}{1!} + \frac{\theta^2}{2!} + \frac{\theta^3}{3!} + \frac{\theta^4}{4!} + \frac{\theta^5}{5!} + \dots$$

Expanding $e^{i\theta}$ formally according to the last series gives:

$$(A III-20) \quad e^{i\theta} = 1 + i \frac{\theta}{1!} - \frac{\theta^2}{2!} - i \frac{\theta^3}{3!} + \frac{\theta^4}{4!} + i \frac{\theta^5}{5!} - \frac{\theta^6}{6!} - i \frac{\theta^7}{7!} + \dots$$

If the $\sin \theta$ series is multiplied by i and added to the $\cos \theta$ series, the result is the series (A III-20); therefore:

$$(A III-21) \quad e^{i\theta} = \cos \theta + i \sin \theta$$

Now

$$(A III-22) \quad e^z = e^{x+iy} = e^x e^{iy} = e^x (\cos y + i \sin y)$$

This is what is meant by the exponential function $w = e^z$. This is an analytic function, as can be shown by applying the test for the Cauchy-Riemann conditions:

$$e^z = e^x \cos y + i e^x \sin y = u + iv$$

where $u = e^x \cos y$ and $v = e^x \sin y$,

$$(A III-23)$$

$$\frac{\partial u}{\partial x} = e^x \cos y; \quad \frac{\partial v}{\partial y} = e^x \cos y; \quad \frac{\partial u}{\partial y} = -e^x \sin y; \quad \frac{\partial v}{\partial x} = e^x \sin y$$

FUNDAMENTALS OF MAPPING

Since $w = f(x, y) + i\phi(x, y)$ and $z = x + iy$, it follows that to each $z_0 = x_0 + iy_0$, there corresponds a $w_0 = f(x_0, y_0) + i\phi(x_0, y_0)$. A new complex plane can be defined, which there is a u -axis, the axis of reals, and a v -axis, the axis of imaginaries, then one can plot w in this plane as a function of u and v ; this is the w -plane. To each u_0 and v_0 there corresponds a w_0 ; then to each z_0 (in the z -plane) there corresponds a w_0 (in the w -plane). If all the z 's along some curve in the z -plane are plotted in the w -plane, the result is a curve in the w -plane which is called the "map" of the curve in the z -plane. If all the points in a given area in the z -plane are plotted in the w -plane, the result is an area in the w -plane. Thus, both curves and areas can be mapped from one plane into the other (See figure A III-5.)

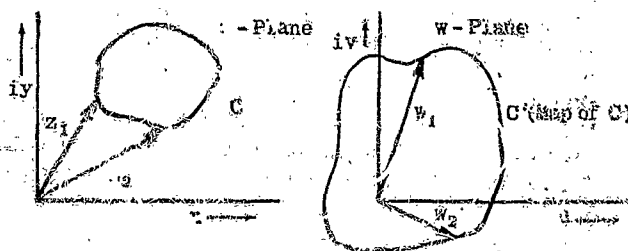


Figure A III-5. Mapping Through $w = f(z)$

Let

$$(A III-24) \quad w(z) = \frac{A_n z^n + A_{n-1} z^{n-1} + \dots + A_0}{B_n z^n + B_{n-1} z^{n-1} + \dots + B_0} = \frac{N(z)}{D(z)}$$

where $N(z)$ and $D(z)$ have no common factor.

It is required to plot this function in the w -plane. In order to do this it is convenient to locate first a curve

points on this plot; for instance, any z_0 which makes the polynomial in the numerator, $n(z)$, vanish gives $w(z_0) = 0$; but the z 's which make the denominator vanish are the roots of the equation

$$(A \text{ III-25}) \quad N(z) = A_m z^m + A_{m-1} z^{m-1} + \dots + A_0 = 0$$

Let those roots be p_1, p_2, \dots, p_m ; then (A III-25) can be factored into

$$(A \text{ III-26}) \quad N(z) = (z - p_1)(z - p_2) \dots (z - p_m)$$

where the p 's are complex roots of $n(z)$. When z assumes each of these values, $N(z)$ vanishes, and so does $w(z)$; therefore, $w(p_1), w(p_2), \dots, w(p_m)$ are all zero. The roots p_1, p_2, \dots, p_m of $N(z)$ are called "zeros" of $w(z)$; they all map into the origin of the w -plane.

Any z_0 which makes the polynomial $D(z)$ vanish causes $w(z)$ to become infinite. These z values are the roots of the equation

$$(A \text{ III-27}) \quad D(z) = B_n z^n + B_{n-1} z^{n-1} + \dots + B_0 = 0$$

Let these roots be s_1, s_2, \dots, s_n ; then

$$(A \text{ III-28}) \quad D(z) = (z - s_1)(z - s_2) \dots (z - s_n)$$

Thus, $w(z)$ becomes infinite for $z = s_1, z = s_2, \dots, z = s_n$; these are called the "poles" of $w(z)$. Since $w(z_1), w(z_2), \dots, w(z_n)$ are infinitely large, they cannot be plotted in the finite portion of the w -plane.

To show them, an "artificial" representation of infinity is introduced; the locus of all infinitely large w 's is shown as a circle, about the origin, of finite radius R ; and it is specified that R approaches infinity.

For all values of z other than the zeros and poles of w it is necessary to substitute a number of z 's into the polynomials of (A III-24) and to solve for the corresponding values of w . If the z 's lie along a closed curve, the plot in the w -plane is also a closed curve; the area inside this curve will plot into an area in the w -plane but, not necessarily inside the w -curve; it can also fall outside the w -curve. The subject of mapping is an important one in the study of servomechanisms. More about it is given in section A IV of this appendix.

INTEGRATION OF FUNCTIONS OF COMPLEX VARIABLES

The integral $\int f(x) dx$ can be thought of as the area under the curve $y = f(x)$. It is the limit of the sum of the areas of the rectangles whose base is Δx and height is $y = f(x)$, as Δx becomes indefinitely small and the number of these rectangles becomes indefinitely large (see figure A III-6).

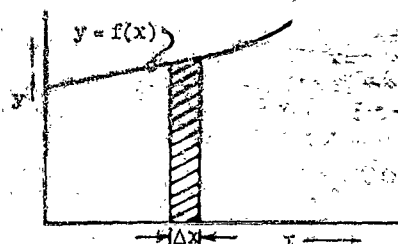


Figure A III-6. Integration as Limit of Summation

Should one form the integral $\int w(z) dz$ and think of it as the limit of the sum of the rectangles whose base is Δz and height $w = f(z)$, the same difficulty arises as in the case of differentiation, namely, in which direction is Δz to be taken. If $\int f(z) dz$ is to have a definite meaning, the direction of Δz has to be prescribed; and this is done by picking out, of the infinite number of z 's in the z -plane, only a set of z 's which lie along some single selected line in the z -plane. For this reason the integral is known as a "line integral"; it is designated by the symbol $\int_C f(z) dz$ (see figure A III-7).

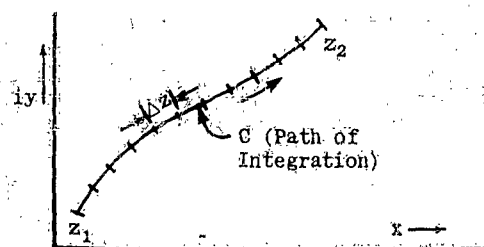


Figure A III-7. Line Integral

If the line has as its end points the values z_1 and z_2 , the integral is a definite one, $\int_{z_1}^{z_2} f(z) dz$. These end points may be anywhere in the z -plane. If the points z_1 and z_2 coincide, i.e., the line is a closed curve, the integral is sometimes designated by the symbol \oint . The line selected in the z -plane is known as the "path of integration" (see figure A III-8).

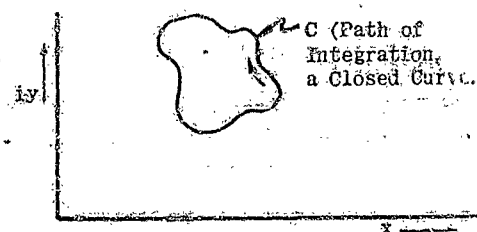


Figure A III-8. Closed Contour

A frequently used path of integration is a circle (see figure A III-9). The equation of such a circle is

$$(A \text{ III-29}) \quad z = z_0 + a e^{i\theta}$$

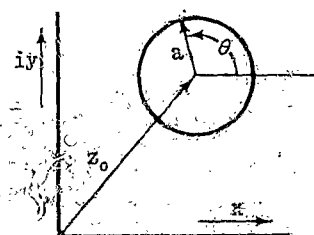


Figure A III-9. Circle $z = z_0 + a e^{i\theta}$

This means that any point on this circle can be located by a vector which is the sum of z_0 and of the vector a whose angle with the positive x -axis is θ .

If the center of the circle is at the origin of the z -plane, the equation (A III-29) simply becomes: (see figure A III-10).

Appendix Section A III

(A III-30)

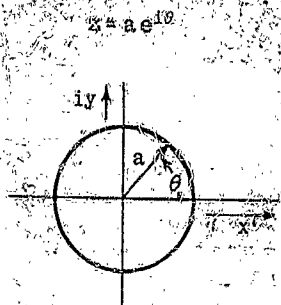


Figure A III-10. Circle $z = ae^{i\theta}$

As an example of a line integration, let $w = z^2$ be integrated along the above circle (figure A III-10), $z = ae^{i\theta}$. For this path $w = a^2e^{2i\theta}$; $dz = ia e^{i\theta} d\theta$; then

$$(A III-31) \oint_C z^2 dz = \int_0^{2\pi} a^2 e^{2i\theta} i a e^{i\theta} d\theta = i a^3 \int_0^{2\pi} e^{i3\theta} d\theta \\ = i a^3 \int_0^{2\pi} \cos 3\theta d\theta - a^3 \int_0^{2\pi} \sin 3\theta d\theta = 0$$

Thus, the integral is zero and is independent of the radius of the circle used as the path of integration.

If the path of integration is the circle shown in figure A III-9, the integral is still zero, as shown below:

$$(A III-32) z = z_0 + ae^{i\theta}; z^2 = z_0^2 + 2z_0 ae^{i\theta} + a^2 e^{2i\theta}; dz = ia e^{i\theta} d\theta;$$

$$\oint_C z^2 dz = i \int_0^{2\pi} z_0^2 a e^{i\theta} d\theta + i \int_0^{2\pi} 2a^2 z_0 e^{i2\theta} d\theta + i \int_0^{2\pi} a^3 e^{i3\theta} d\theta = \\ = i z_0^2 a \int_0^{2\pi} [\cos \theta + i \sin \theta] d\theta + i 2 a^2 z_0 \int_0^{2\pi} [\cos 2\theta + i \sin 2\theta] d\theta + i a^3 \int_0^{2\pi} [\cos 3\theta + i \sin 3\theta] d\theta = 0$$

In this case, too, the result is independent of the radius; nor does the position of the center of the circle matter.

As another example, let the same function $w = z^2$ be integrated along a path which is a square in the z -plane (see figure A III-11(a)).

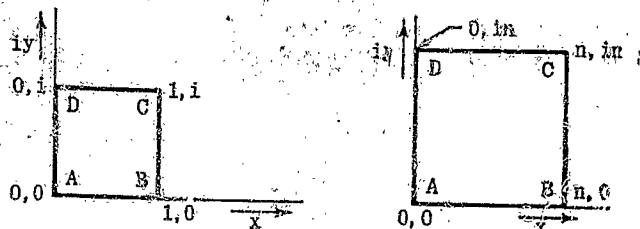


Figure A III-11. Rectangles as Paths of Integration

It is seen that along AB $z = x$; $dz = dx$, since y does not change; along BC $z = 1 + iy$; $dz = i dy$, since x does not change; along DC $z = x + i$; $dz = dx$; along AD $z = iy$; $dz = i dy$.

Integrating along the sides of the square, counterclockwise, gives for the integrals along AB, BC, CD and DA respectively;

$$(A III-33) \oint_C z^2 dz = \int_0^1 x^2 dx + \int_0^1 (1 + iy)^2 i dy + \int_1^0 (x + i)^2 dx - \int_1^0 (iy)^2 i dy \\ = \int_0^1 x^2 dx + i \int_0^1 (1 + 2iy - y^2) dy - \int_1^0 (x^2 + 2ix - 1) dx \\ = \int_0^1 x^2 dx + i \int_0^1 2x dx + i \int_0^1 dy - \int_0^1 2y dy - \int_1^0 y^2 dy - \int_1^0 x^2 dx \\ = \int_0^1 x^2 dx + \int_0^1 2x dx + i \int_0^1 y^2 dy = 1 - 1 + i = 0$$

It is seen that for $w = z^2$, the value is zero, whether C is a large circle, a small one or even a square. Had the square been n times as large (see figure A III-11 (b)) the limits would become zero and n ; along BC $z = n + iy$ along CD $z = x + in$; the integral now would be:

(A III-34)

$$\oint_C z^2 dz = \int_0^n x^2 dx + \int_0^n (n + iy)^2 i dy + \int_n^0 (x + in)^2 dx + \int_n^0 (iy)^2 i dy \\ = \int_0^n x^2 dx + in^2 \int_0^n dy + in \int_0^n 2y dy + i \int_0^n y^2 dy \\ + \int_n^0 x^2 dx + in^2 \int_n^0 dy + in \int_n^0 2y dy + i \int_n^0 y^2 dy \\ = \int_0^n x^2 dx + in^2 \int_0^n dy - n \int_0^n 2y dy - i \int_0^n y^2 dy - \int_n^0 x^2 dx \\ - in^2 \int_n^0 dy - in \int_n^0 2y dy - i \int_n^0 y^2 dy \\ = in^3 - n^3 - in^3 + n^3 = 0$$

It is clear that $\oint_C z^2 dz$ is zero for any finite size of square. By trying other paths of integrations it can be found that $\oint_C z^2 dz$ is zero for any path. This is true because z^2 is analytic at any finite point in the z -plane. And it can be shown that for any function $f(z)$ which is analytic at any point in the z -plane the integral along any closed curve is zero.

(A III-35)

$$\oint_C f(z) dz = 0$$

for any closed path if $f(z)$ is analytic anywhere in the z -plane.

This property of analytic functions which is given here without proof, but only made plausible by a few examples, is the content of the Cauchy-Goursat theorem.

Such a property of a function is without precedent in the realm of real numbers and might, therefore, prove somewhat puzzling. To make it easier to accept it, one may think of some physical models in which an integral of a function exhibits similar characteristics. For example, let a tension spring be attached at the point (0,0) (see figure A III-12).

The free end, A, is made to follow the curve C. At any point z along the curve the extension of the spring is $|z|$. The potential energy E is proportional to the extension. Thus E is a function of z. The net increase in energy at any point is the summation of the changes in energy up to that point. The total change in energy as A returns to its initial point is the integral of the change dE along the curve. And, since the spring ends up in its initial condition, there is no

Therefore, the integral is zero. For a closed path, the change in potential energy is zero. The net change in potential energy at any point along the path is the integral of the changes up to that point. If the man returns to his starting point, the change in his potential energy is zero. This is another integral which is zero, regardless of the path of integration.

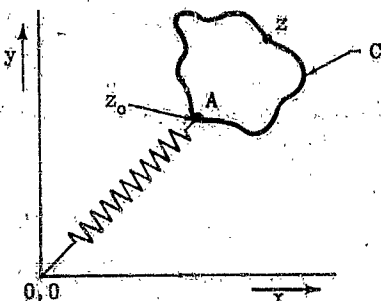


Figure A III-12. Tension Spring as Model of Analytic Function

As a final example, let a torsion spring, (in the shape of a rod) have one end attached to the point (0,0) (see figure A III-13).

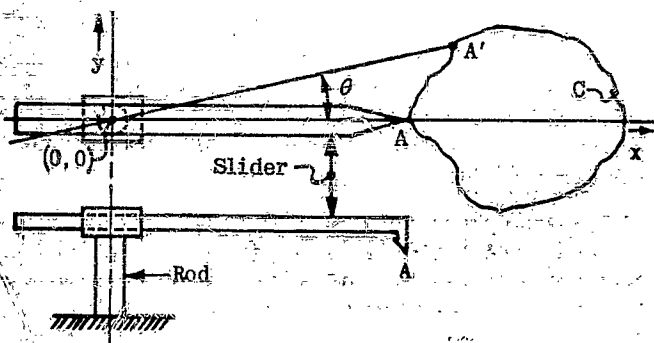


Figure A III-13. Torsion Spring as Model of Analytic Function

Let the point A of the slider move along a closed curve, such as the one in figure A III-13. The torsional energy of the rod is proportional to the angle θ ; as the point A comes back to its original position the net energy in the torsion spring (i.e., in the rod) is zero and this is true for any path not surrounding the origin, i.e., the point where the rod is attached.

Turning to the function $w=1/z$, it is seen that it becomes infinite at $z=0$ (i.e., it has a pole at $z=0$). This function is analytic everywhere in the z -plane except at the origin. It is of interest to integrate this function along a circle having its center at the origin (see figure A III-10).

Thus along C $z=ae^{i\theta}$; $dz=iae^{i\theta}d\theta$; and

$$(A III-36) \quad \int_C \frac{1}{z} dz = \int_0^{2\pi} \frac{iae^{i\theta} d\theta}{ae^{i\theta}} = i \int_0^{2\pi} d\theta = 2\pi i$$

It is seen that the integral in this case is not zero but $2\pi i$. For $w=1/z^n$,

$$(A III-37) \quad \int_C \frac{1}{z^n} dz = \int_0^{2\pi} \frac{iae^{i\theta} d\theta}{a^n e^{in\theta}} = \frac{1}{a^{n-1}} \int_0^{2\pi} e^{-i(n-1)\theta} d\theta = \frac{1}{a^{n-1}} \int_0^{2\pi} [\cos(m-1)\theta - i \sin(m-1)\theta] d\theta = 0$$

for all $m \neq 1$.

Thus, $\int_C z^n dz$ is zero for any n , positive or negative, except $n=-1$, for which its value is $2\pi i$.

If the path of integration is an arc of the above circle with the angle of the end points z_1 and z_2 being θ_1 and θ_2 (see figure A III-14), the integral becomes:

$$(A III-38) \quad \int_{z_1}^{z_2} \frac{1}{z} dz = i \int_{\theta_1}^{\theta_2} d\theta = i(\theta_2 - \theta_1)$$

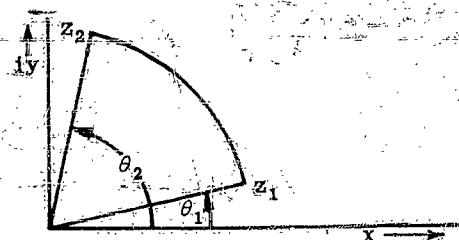


Figure A III-14. Integration Along Circular Arc

if the integration proceeds from z_1 to z_2 , and $-i(\theta_2 - \theta_1) = i(\theta_1 - \theta_2)$ if it proceeds from z_2 to z_1 .

The integral of $1/z^n$ (where $n \neq 1$) over the same arc gives

$$(A III-39) \quad \frac{1}{a^{n-1}} \int_{\theta_1}^{\theta_2} \{\cos[(n-1)\theta] - i \sin[(n-1)\theta]\} d\theta$$

As the radius a increases, the value of (A III-39) decreases and as a approaches infinity, it approaches zero.

The integral of z^n over this arc is:

$$(A III-40)$$

$$\int_{z_1}^{z_2} z^n dz = \int_{\theta_1}^{\theta_2} a^n e^{in\theta} \cdot iae^{i\theta} d\theta = ia^{n+1} \int_{\theta_1}^{\theta_2} e^{i(n+1)\theta} d\theta = ia^{n+1} \int_{\theta_1}^{\theta_2} \{\cos[(n+1)\theta] + i \sin[(n+1)\theta]\} d\theta$$

This integral varies as a^{n+1} and as a approaches zero, it also approaches zero.

To summarize, the integral of z^n over an arc of radius a :

$$(A III-41) \quad a) \text{ for } n < -2 \text{ and } a \text{ very large, } \int_{z_1}^{z_2} z^n dz \rightarrow 0$$

$$b) \text{ for } n = -1 \text{ and any } a, \int_{z_1}^{z_2} z^n dz = i(\theta_2 - \theta_1)$$

$$c) \text{ for } n > 0 \text{ and } a \text{ very small } \int_{z_1}^{z_2} z^n dz \rightarrow 0$$

In taking the integral along a simple closed curve (i.e.,

Appendix Section A IV

one that does not cross itself,) the counterclockwise direction is taken as positive.

It has been shown that the value of a line integral of a function depends not only on the function itself but also on the path of integration. It will now be shown that when the path of integration is a closed curve and, furthermore, the function is analytic on the curve and at any point in the region enclosed by this curve, except at the point z_0 (that means, z_0 is a pole of this function), then the value of this line integral is a known function of this point z_0 .

This can be done as follows: The function having a pole at $z = z_0$ is of the form $w(z) = f(z)/(z - z_0)$ (see figure A III-15), where $f(z)$ does not contain the factor $(z - z_0)$.

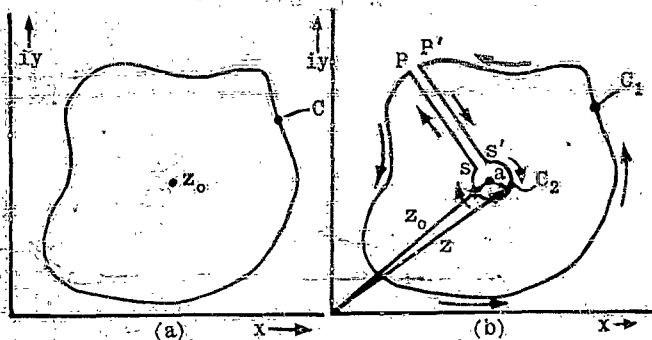


Figure A III-15. Contour of Integration Enclosing a Pole

Because of the presence of this pole, the line integral around C is not zero. But, one can change the contour C in such a manner, as to avoid the pole; (see figure A III-15(b)). Now the function is analytic on the new contour and at all points inside the enclosed region. Therefore, the Cauchy-Goursat theorem applies, and carrying out the integration as shown in the figure, one can write:

(A III-42)

$$\oint_C w(z) dz = \int_C w(z) dz + \int_P^{P'} w(z) dz + \int_{C_2}^{P'} w(z) dz - \int_C w(z) dz$$

(The negative sign is used because C_2 is traversed in a direction opposite to C_1). Since PS can be made indefinitely close to $S'P'$, and since the two lines are traversed in opposite directions, the integral over this part of the path cancels out. Noting that the integral around the circle C_2 is taken in a direction opposite to that of C_1 , there remains:

$$(A III-43) \quad \int_{C_1} w(z) dz = \int_{C_2} w(z) dz$$

The equation of the circle C_2 is

$$(A III-44) \quad z = z_0 + ae^{j\theta}$$

Also: $dz = ja e^{j\theta} d\theta$ and $z - z_0 = ae^{j\theta}$. By substituting these values into the right half of equation (A III-43),

(A III-45)

$$\int_{C_2} \frac{f(z) dz}{z - z_0} = \int_0^{2\pi} \frac{f(z_0 + ae^{j\theta}) ja e^{j\theta} d\theta}{ae^{j\theta}} = j \int_0^{2\pi} f(z_0 + ae^{j\theta}) d\theta$$

The radius a can be made arbitrarily small; in the limit it approaches zero. Then equation (A III-45) becomes:

$$(A III-46) \quad \int_{C_2} \frac{f(z) dz}{z - z_0} = 2\pi j f(z_0);$$

and, finally, equation (A III-43) becomes:

$$(A III-47) \quad \int_{C_1} \frac{f(z) dz}{z - z_0} = 2\pi j f(z_0);$$

This relation is known as the Cauchy integral formula. It states that if a function has a pole at some point z_0 (i. e., it is of the form $[f(z)]/(z - z_0)$), the integral of this function around any closed curve surrounding the pole is equal to $2\pi j f(z_0)$ where $f(z_0)$ is the value of $f(z)$ at the pole.

As an example of the application of the Cauchy integral formula, let $w(z) = \sin z/(z - z_0)$ (which shows that $w(z)$ has a pole at $z = z_0$). Then the integral of $w(z)$ around a path which encircles the point z_0 is:

$$\oint w(z) dz = \oint \frac{\sin z}{z - z_0} dz = 2\pi j \sin z_0$$

$$= 2\pi j [\sin x_0 \cosh y_0 + j \cos x_0 \sinh y_0]$$

SECTION A IV - MAPPING

(a) INTRODUCTION

One of the methods for investigating the stability of a servomechanism is based on the use of the so-called "Nyquist criterion." In order to understand the derivation of this criterion it is necessary to know something about the operation of "mapping" of a given contour from the z -plane into the w -plane. Furthermore, before a discussion of mapping may be begun certain fundamental concepts regarding plots in a complex plane must be understood. For these reasons the presentation in this chapter will be as follows:

1. Fundamental concepts—methods for designating points, lines and areas in a complex plane.
2. Mapping.
3. Nyquist criterion.

(b) METHODS OF DESIGNATING POINTS, LINES AND AREAS IN THE Z -PLANE

It is important to become familiar with the ways used to designate points, lines and areas in the z -plane.

In the theory of complex numbers it is shown that a fixed point in the z -plane is designated by $z_0 = x_0 + jy_0$ (see figure A IV-1).

Straight lines when unlimited, are designated as follows: The y axis is $x = 0$; the x axis is $y = 0$; a line parallel to the y axis and at a distance C from it is $x = C$; when $C > 0$, the line is in the right half plane (see figure A IV-2). The designation of certain other lines is shown in figure A IV-2.

The half plane lying to the right of the line $x = C$ is $x > C$ (see figure A IV-3).

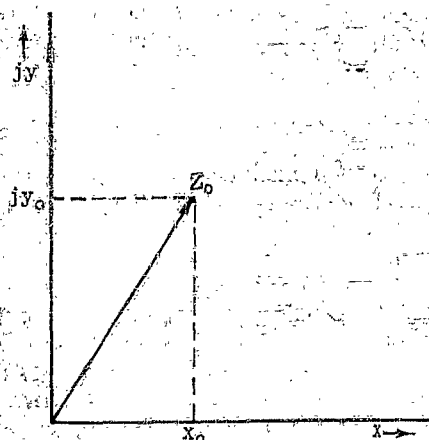


Figure A IV-1. Point z_0

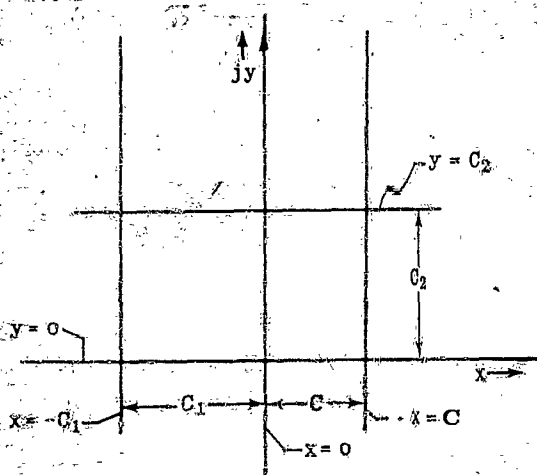


Figure A IV-2. Straight lines

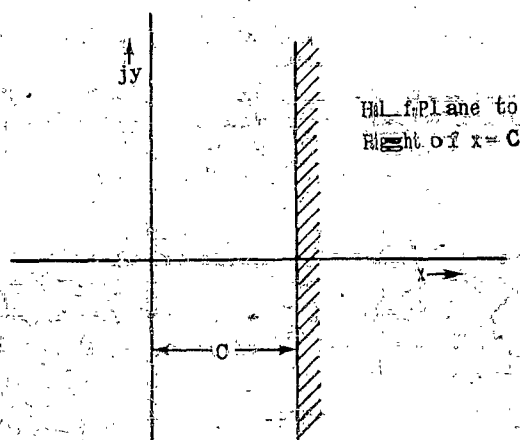


Figure A IV-3. Half-Plane

It follows that $x > 0$ is the entire right half-plane. (Note the crosshatching used to designate an area.) Similarly $y > 0$ is the entire half-plane lying above the x axis.

Straight semi-infinite lines starting at the origin are designated by their phase angle, θ (see figure A IV-4).

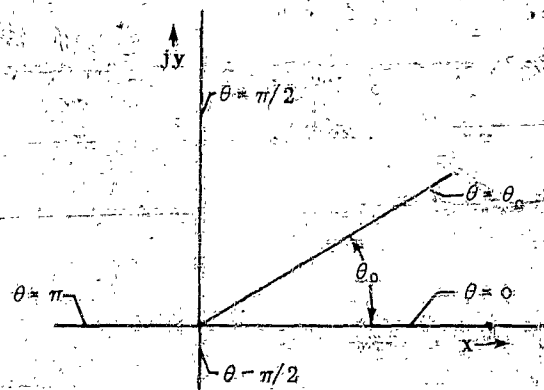


Figure A IV-4. Radial Straight Line

Thus, the positive half of the x -axis is $\theta = 0$; the negative half is $\theta = \pi$.

Partly bounded areas are designated as shown in figure A IV-5.

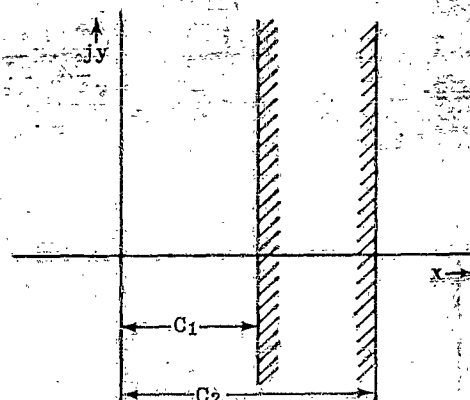


Figure A IV-5. Infinite Strip

A vertical strip is $C_1 < x < C_2$; a horizontal strip is $C_3 < y < C_4$. Using the phase angle, the upper half-plane can be designated by $0 < \theta < \pi$; and the right half-plane by $-\pi/2 < \theta < +\pi/2$; if the boundary lines are included, one writes $-\pi/2 \leq \theta \leq +\pi/2$. The first quadrant is either $0 < \theta < \pi/2$ or $x > 0, y > 0$; the second quadrant is $\pi/2 < \theta < \pi$ or $x < 0, y > 0$.

A rectangle can be defined by its boundaries: $x_1 = C_1$,

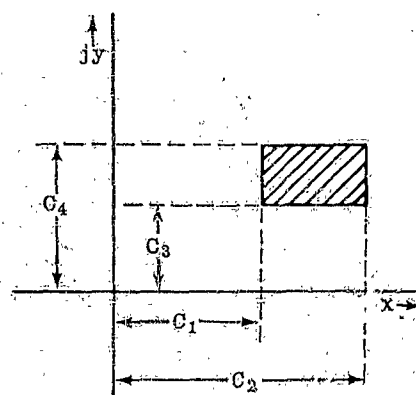


Figure A IV-6. Rectangle

Appendix Section AIV

$x_2 = C_2$; $y_1 = C_3$; $y_2 = C_4$. The area inside this rectangle is: $C_1 < x < C_2$; $C_3 < y < C_4$. A circle of unit radius is $r = 1$ (see figures A IV-6 and A IV-7).

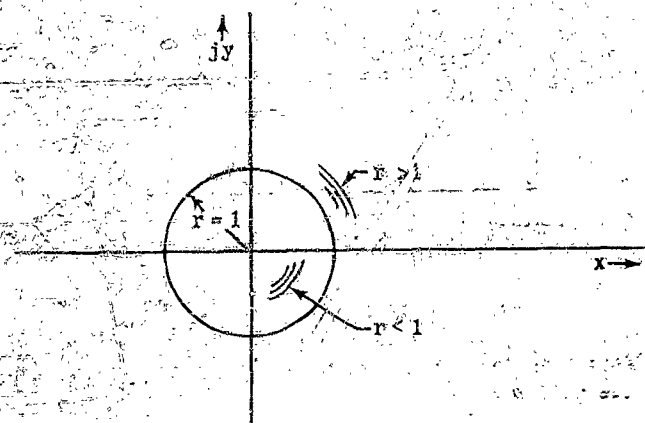


Figure A IV-7. Circle at Origin

The area inside the unit circle is $r < 1$; if the boundary is included, one writes: $r \leq 1$; the entire area outside the unit circle is $r > 1$; another way of designating the unit central circle is: $z = e^{j\theta}$. A circle of radius a and with its center at $z = z_0$ is $z = z_0 + ae^{j\theta}$ (see figure A IV-8).

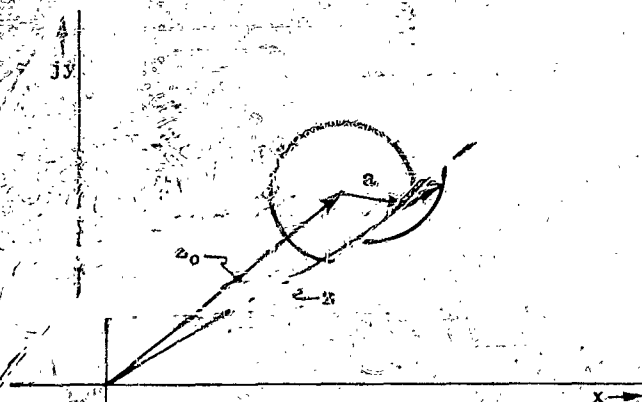


Figure A IV-8. General Circle

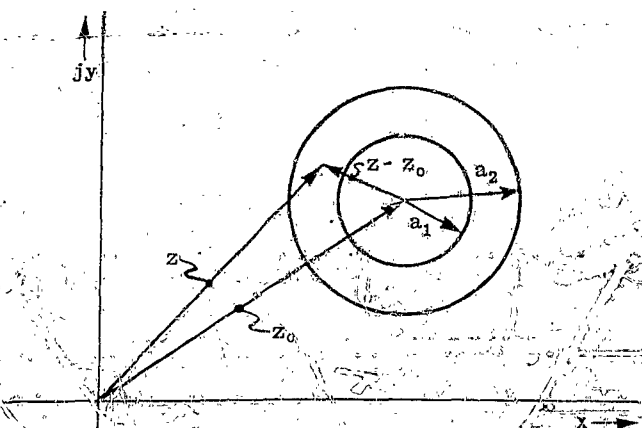


Figure A IV-9. Annulus

The area inside the circle is $z = z_0 + ae^{j\theta}$.

A ring-shaped area bounded by the concentric circles of radii a_1 and a_2 and having their center at z_0 is: (see figure A IV-9)

$$a_1 < |z - z_0| < a_2$$

The above examples should be sufficient to enable one to identify all the designations to be used later.

(c) MAPPING

The problem of mapping is, essentially, as follows: The problem of mapping is to draw the corresponding contour in the w -plane. If the function $w(z)$ is a simple one, this mapping can be done by transforming the contour right on the z -plane. But, generally, a separate w -plane must be drawn.

MAPPING A CONTour THROUGH THE FUNCTION $w = z + C_0$ (see figure A IV-10).

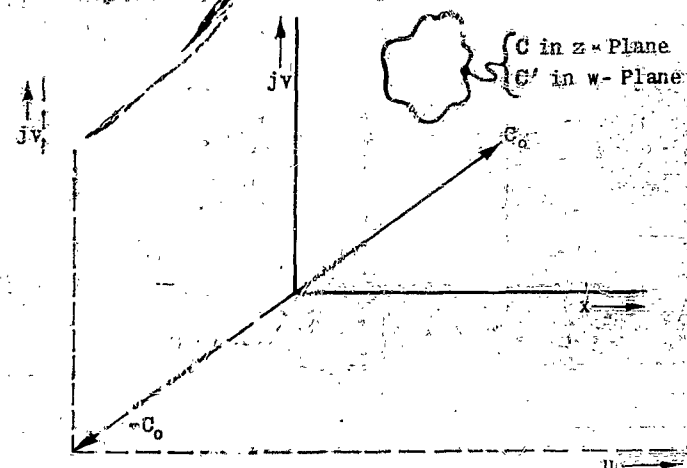


Figure A IV-10. Mapping Through $w = z + C_0$

Let there be an irregular contour C in the z -plane. To each point on this contour corresponds a point $w = z + C_0$, where C_0 is a complex number. It follows that w transforms C merely by shifting it in its entirety through the distance $|C_0|$, in the direction of the angle of C_0 . Rather than draw the contour C in its new position it is simpler to shift the axes to a point $-C_0$. Thus, with respect to the uv axes, that is, in the w -plane the map of C is of the same shape and orientation as C , but shifted through C_0 . This map or "image" of C is designated by C' .

Analytically, this transformation is arrived at as follows:

$$(A IV-1) \quad \begin{cases} C_0 = x_0 + jy_0; & z = x + jy; & w = z + C_0 = (x + x_0) + j(y + y_0) \\ w = u + jv; & u = x + x_0; & v = y + y_0 \end{cases}$$

This shows that every point of the contour C is shifted horizontally through x_0 and vertically through y_0 , hence, through the distance $|C_0|$, in the proper direction, as was done graphically. This simple transfor-

mation helps to explain a theorem which is used later on.

Let $w = z - C_0$, and let C_0 be a point inside the contour C in the z -plane (see figure A IV-11).

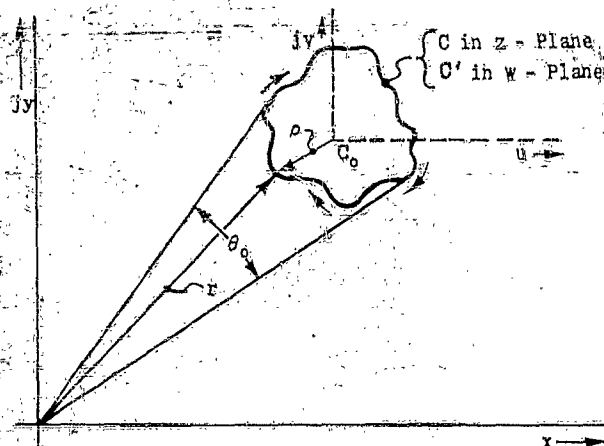


Figure A IV-11. Mapping (Translation)

Then, rather than shift the contour C through $-C_0$, the uv axes are shifted by $+C_0$, i.e., the origin of the uv coordinates is at point C_0 . The contour C is now also the contour C' in the w -plane. (Note that plotting a contour through the function $w = z - C_0$ amounts to putting the origin of the uv axes at the point C_0 .)

Now, let a point z traverse the contour C in a clockwise direction; using $z = re^{j\theta}$, it is apparent that, as z travels completely around the contour C , the vector r swings back and forth through an included angle θ_0 but ends up in its initial position so that the net angle traversed is zero. Using $w = \rho e^{j\phi}$, it is likewise apparent that the vector ρ rotates through 360° about the origin as the point w travels around C' ; w travels around C' in the same direction as z travels around C .

A theorem can now be stated as follows: When a closed contour C is mapped through the function $w = z - C_0$, where C_0 lies inside the contour, then, as z travels completely around C , w travels once around the origin and in the same direction.

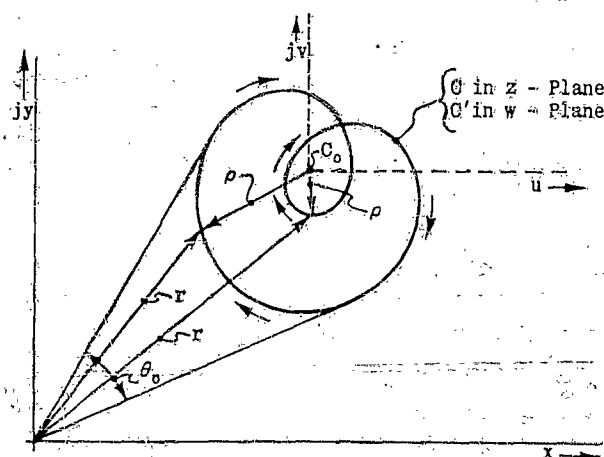


Figure A IV-12. Mapping (Multiple Encirclement)

This theorem can be extended to the case where the contour C has a loop and thereby encircles the point C_0 twice (see figure A IV-12). It is seen that as z travels completely around C , w travels twice around the origin. In general, when C encircles C_0 n times, w goes around the origin n times and ϕ rotates through $2\pi n$ radians.

If the point C_0 is outside the contour C , then w does not go around the origin but swings back and forth through some angle θ_0 with net angle traversed equal to zero (see figure A IV-13).

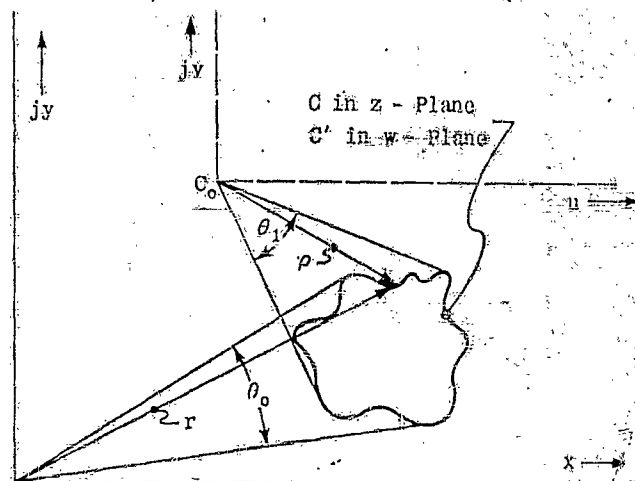


Figure A IV-13. Mapping

MAPPING OF A CONTOUR THROUGH THE FUNCTION $w = C_0 z$. To map a contour C in the w -plane through the function $w = C_0 z$, where C_0 is a complex number, it is best to use the polar form for z and w ; thus:

$$(A IV-2) \quad \begin{cases} z = re^{j\theta}; & w = \rho e^{j\phi}; & C_0 = r_0 e^{j\theta_0} \\ w = C_0 z = r_0 r e^{j(\theta_0 + \theta)}; & \therefore \rho = r_0 r; & \phi = \theta_0 + \theta \end{cases}$$

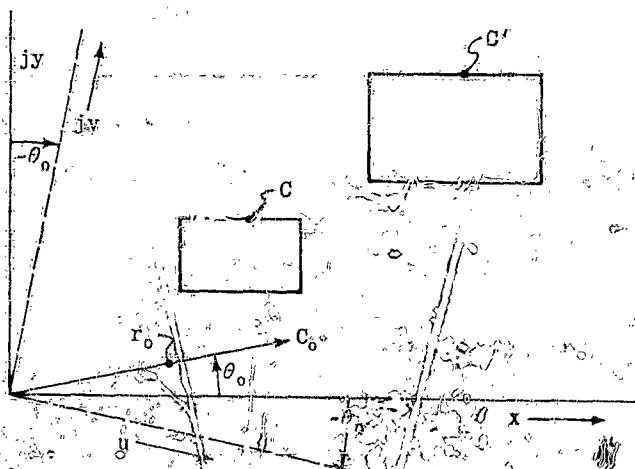


Figure A IV-14. Mapping Through $w = C_0 z$

Thus, the transformation consists in enlarging all vectors r by a factor r_0 and then rotating these enlarged vectors through the angle θ_0 (see figure A IV-14).

Appendix
Section A IV

(In this figure r_0 was taken to be about 2 units. Instead of rotating the contour C , the u and v axes were rotated through $-\theta_0$.)

The result of this transformation is geometrically similar to C , is a contour C' which is removed from the origin C_0 , but enlarged (if $r_0 > 1$), relative to the uv axes and with its orientation relative to its original orientation by a constant angle from

where C_1 and C_2 are complex numbers. Calling

$$(A IV-4) \quad w_1 = z - C_1 = \rho_1 e^{j\phi_1}; \text{ and } w_2 = z - C_2 = \rho_2 e^{j\phi_2}$$

one can write (see figure A IV-15):

$$(A IV-5) \quad \begin{cases} w = \rho e^{j\phi} = w_1 w_2 = \rho_1 \rho_2 e^{j(\phi_1 + \phi_2)} \\ \therefore \rho = \rho_1 \rho_2; \quad \phi = \phi_1 + \phi_2 \end{cases}$$

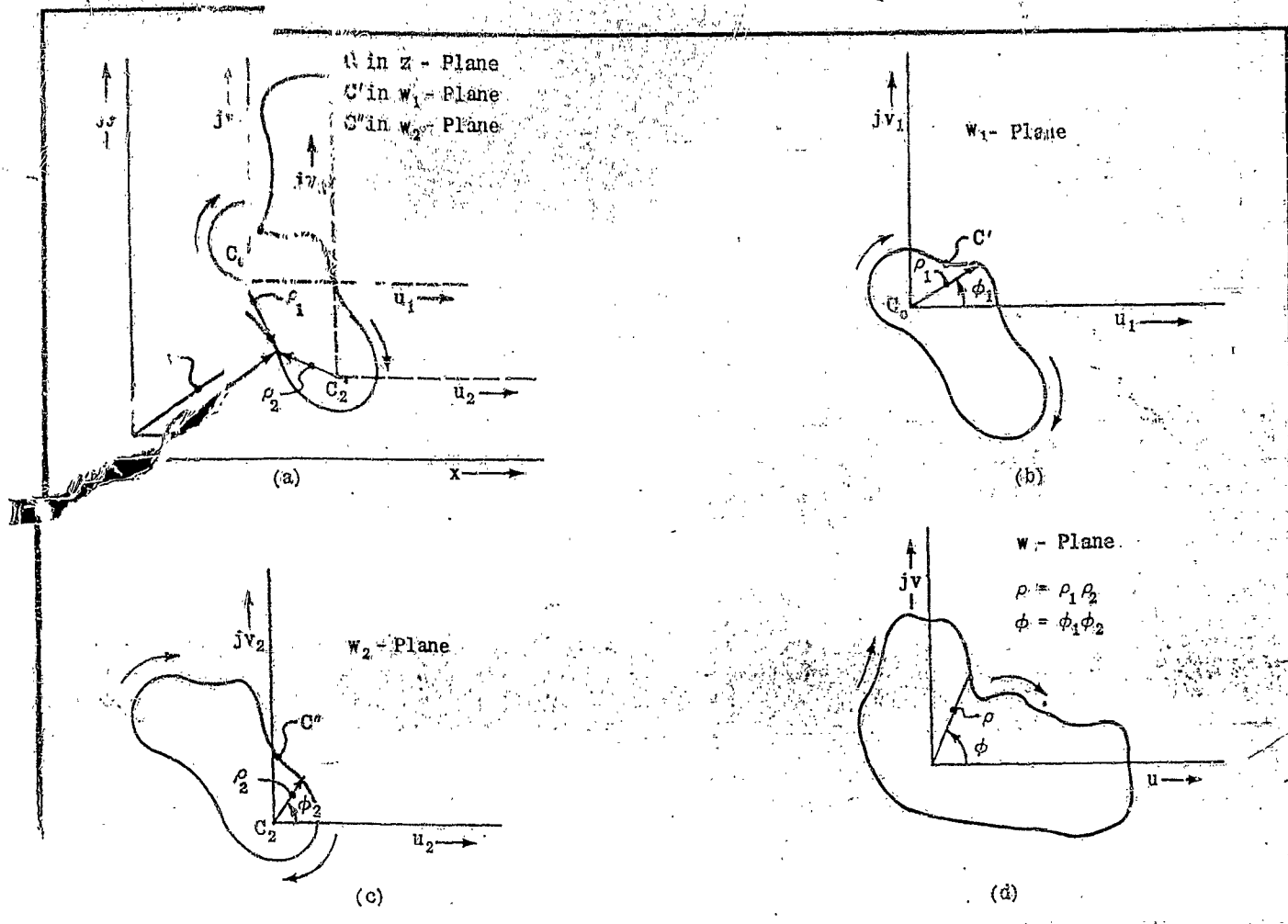


Figure A IV-15. Mapping Through (A IV-5).

the above two transformations can be combined into

A moving to this transformation a contour C is first complicated by C_0 , giving $w' = C_0 z$, then the resulting contour is shifted by C_1 , giving $w = w' + C_1$. Graphically this can be done by plotting C into C' , as shown in Figure A IV-14 and then shifting the origin by $-C_1$, as in Figure A IV-10.

In all the transformation described so far C' remains geometrically similar to C , though its orientation relative to the uv axes may be different from that of C relative to the xy axes.

The following transformation can be further extended as follows. Let

$$w = (z - C_1)(z - C_2)$$

Examination of w_1 shows that if C_1 lies inside the contour C , ρ_1 traveling along C' in the w_1 -plane rotates through 2π radians around the w_1 origin as z goes completely around C . Similarly, ρ_2 traveling along C'' in the w_2 -plane rotates through 2π radians about the w_2 origin as z goes completely around C .

Consequently, when ϕ_1 and ϕ_2 each goes through 2π radians, ϕ which is $\phi_1 + \phi_2$ goes through 4π radians. The theorem is now extended to the case when C contains two points, C_1 and C_2 , and the contour is mapped through $w = (z - C_1)(z - C_2)$. Then, as z goes once around the contour, w goes $2\pi \times 2$ times around the origin, and in the same direction as z (i.e., clockwise, if z travels clockwise).

If both C_1 and C_2 are outside the contour C , w does not go around the origin but rotates through an angle

$\phi = \phi_1 + \phi_2$. If one of the points lies inside C , w goes one complete time around the origin. Similarly, for m points inside the contour C , w goes through $2\pi m$ radians around the origin, i. e., it encircles the origin m complete times, rotating in the same direction as z .

The extended theorem can now be stated: When a closed contour is mapped through a function $w = (z - C_1)(z - C_2) \dots (z - C_m)$ where all C_i 's lie inside the contour, then as z travels completely around C , w travels around the origin m times and in the same direction.

There is no restriction on the location of C_1 and C_2 , as long as they are inside the contour C . Thus, C_1 may coincide with C , giving $w = (z - C_1)^2$; the rule still remains the same.

MAPPING A CONTOUR THROUGH THE FUNCTION $w = 1/z$. It will be shown that when the contour G is mapped through this function, the resulting contour C' is no longer geometrically similar to C but is generally, distorted into quite a different shape. This transformation still transforms circles into circles, but straight lines become either circles or straight lines, depending on their position in the z -plane.

Using the polar form, one writes:

$$(A \text{ IV-6}) \quad \begin{aligned} z &= re^{j\theta}; \quad w = \rho e^{j\phi} = \frac{1}{z} = \frac{1}{r} e^{-j\theta} \\ \therefore \rho &= \frac{1}{r}; \quad \phi = -\theta \end{aligned}$$

These relations explain the reason for the distortion of the contour: each vector r turns into its own reciprocal, $1/r$; thus, points far away from the origin come quite near, and vice versa. Moreover, each angle θ turns into its negative, $-\theta$. Thus, the contour is reflected across the x -axis, i. e., points of C lying in the upper half-plane become points in C' lying in the lower half-plane (see figure A IV-16).

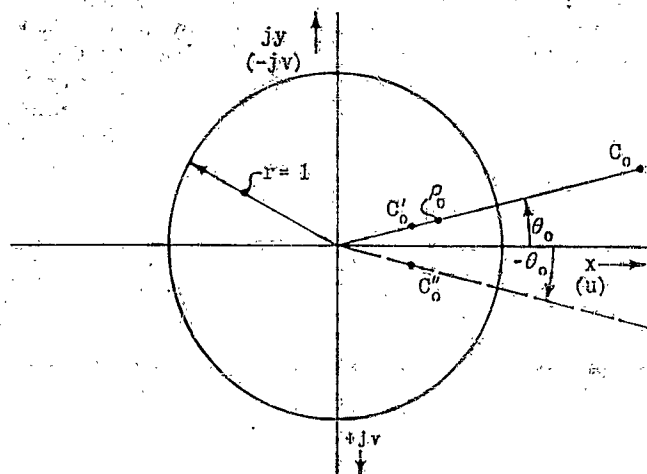


Figure A IV-16. Mapping; $w = 1/z$

Figure A IV-16 shows the mapping of a single point, C_0 into the w -plane. It also shows the unit circle, $r = 1$. r_0 is shown to be about 2 units; ρ_0 is $1/r = 1/2$ unit. The final image, C'_0 , is the conjugate of C_0 .

Figure A IV-17 is the image of the circle $r = 2$.

This image is a circle of radius $\rho = 1/2$; a point z moving along the contour C from A to B moves clockwise; the corresponding point w moving from A' to B' moves counterclockwise. It is important to note that this is the opposite of the behavior of the transformations discussed before.

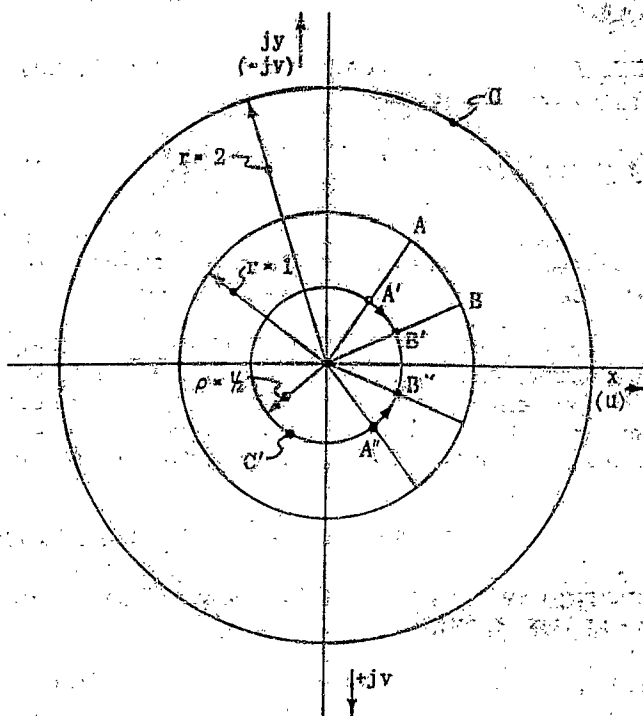


Figure A IV-17. Mapping; $w = 1/z$

The "encirclement" theorem can now be extended once more. Let

$$(A \text{ IV-7}) \quad \begin{cases} w = \frac{1}{z - C_0} = \rho e^{j\phi}; \quad w_1 - z = \rho e^{j\phi_1} \\ \therefore w = \frac{1}{w_1} = \frac{1}{\rho_1} e^{-j\phi_1} \end{cases}$$

Examination of $w = 1/(z - C_0)$ (C_0 inside contour C) shows that, as z moves around C in a clockwise direction once, w_1 circles the origin once, also in a clockwise direction. And since $w = 1/\rho_1 e^{-j\phi_1}$, it follows that, as w_1 circles the origin once in a clockwise direction, w circles its origin once in a counterclockwise direction. By analogy one may infer that for

$$(A \text{ IV-8}) \quad w = \frac{1}{(z - C_1)(z - C_2)}$$

when C_1 and C_2 lie inside the contour C then, as z travels around C once in a clockwise direction, w rotates through $2\pi \times 2$ radians around its origin, but in the counterclockwise direction. Once more, if C_2 coincides with C_1 , $w = 1/(z - C_1)^2$; the rule still remains the same.

MAPPING THEOREM UNDERLYING NYQUIST CRITERION OF STABILITY. The above theory is sufficient to prove an important theorem used in servomechanism work.

Appendix Section A V

It will be recalled (see section A III, Functions of Complex Variable) that a transfer function can be of the form

$$(A\ IV-9) \quad w(z) = \frac{N(z)}{D(z)} = \frac{(z-a_1)(z-a_2)\dots(z-a_n)}{(z-\alpha_1)(z-\alpha_2)\dots(z-\alpha_n)}$$

a_1, a_2, \dots, a_n are the zeros of $w(z)$; $\alpha_1, \alpha_2, \dots, \alpha_n$ are its poles.

If there are multiple roots of $N(z)$ or of $D(s)$, this expression may become:

$$(A\ IV-10) \quad w(z) = \frac{(z-a_1)^K(z-a_2)\dots(z-a_n)}{(z-\alpha_1)^F(z-\alpha_2)\dots(z-\alpha_n)}$$

Let there be a contour C in the z -plane and let all the zeros and poles of $w(z)$ lie inside this contour; let C be mapped into the w -plane through the function $w(z)$ shown in (A IV-10). It will at once be seen, as a consequence of the above "encirclement" theorem, that as z travels around C once, in a clockwise direction, the point w traveling on its image C' , will circle the origin of the w -plane in the same direction as many times as there are roots a , (i.e., zeros) and in the

opposite direction as many times as there are roots α , (i.e., poles). If there are Z zeros and P poles lying inside C , the n , calling the number of encirclements of the origin of the w -plane, in the same direction as that of Z , N , this theorem states that:

$$(A\ IV-11) \quad N = Z - P$$

Defining the counterclockwise direction as positive, (A IV-11) can be expressed in words as follows: "If the contour C encircles Z zeros and P poles in a positive sense, the contour C' encircles the origin $N = Z - P$ times in a positive sense."

In servo work $w(z)$ usually is of the form $w(z) = 1 + Y(z)$. It is more convenient to work with $Y(z)$ and thus the $Y(z)$ -plane. Since $Y(z) = w(z) - 1$, any point in the $w(z)$ -plane maps into the $Y(z)$ -plane shifted to the left by one unit in the $Y(z)$ -plane. Thus the origin of the $w(z)$ -plane becomes the -1 point of the $Y(z)$ -plane and the mapping theorem may be restated.

"If the contour C encircles Z zeros and P poles in a positive sense, the contour C' encircles the -1 point $N = Z - P$ times in a positive sense."

SECTION A V — FACTORING POLYNOMIALS BY SERVO ANALYSIS METHODS

SECTION AV — FACTORING POLYNOMIALS BY SERVO ANALYSIS METHODS

Any of the servo analysis methods developed in Chapter III can be used to find the roots of equations to varying degrees of approximation. This is accomplished by rearranging polynomials to obtain a succession of equations of the form $1 + f(x) = 0$ and applying the special methods developed in that chapter for finding roots of this type of equation. A method is explained here in detail using the root locus method only as a matter of convenience.

Any rational polynomial with constant coefficients can be factored by the root locus method. In order to do this the equation is rearranged as in the following steps:

$$\begin{aligned} (A\ V-1) \quad & x^4 + Ax^3 + Bx^2 + Cx + D = 0 \\ & (x+A)x^3 + Bx^2 + Cx + D = 0 \\ & [(x+A)x + B]x^2 + Cx + D = 0 \\ & \{[(x+A)x + B](x+C) + D\}x = 0 \end{aligned}$$

Next, the expression in the inner brackets $[\]$ is solved as follows:

$$\begin{aligned} (A\ V-2) \quad & (x+A)x + B = 0 \\ & \frac{(x+A)x}{B} + 1 = 0 \\ & \frac{1}{B}(x+A)x = -1 \end{aligned}$$

The last expression may be represented by the complex number $re^{j\theta}$. This procedure is the same as was used in plotting the root locus of $Y(s) = -1$. Thus, using $x = 0$ and $x = -A$ as the zeros of the last equation in

(A V-2), one can plot a locus of points for which $\angle \phi_N = 180^\circ$ (see figure A V-1). Next, on this locus one finds two points, $-a + j\omega$ and $-b + j\omega$ which satisfy the condition $|Kx_N| = 1$. In this case $K = 1/B$; there are just two r_N 's: $-a$ and $-b$.

(It will be seen that the expression $1/B(x+A)x = -1$ is different from those encountered in the development of the root locus method in section III-4 in that it has no denominator. The situation was not covered in section III-4 which deals only with systems in which the order of the numerator is lower than that of the denominator. However, this property, as well as some others dependent upon the poles of $Y(s)$ pertains only to transfer functions of physical systems. But the root locus method is not limited to such systems; it can be applied to purely mathematical equations, i.e., ones which need not be tied-up with the dynamics of physical systems. The form of the root locus of these equations will differ in some respects from those shown previously.)

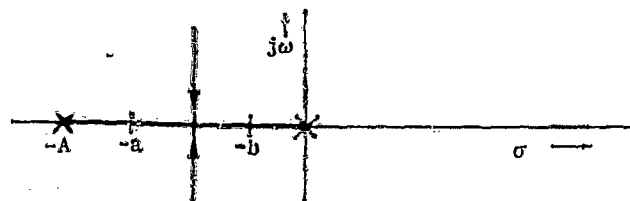


Figure A V-1. $(x+A)x + B = 0 = (x+a)(x+b)$

Note that the locus comes in from infinity and terminates at the zeros as B decreases from infinity. Having determined roots $-a$ and $-b$, the next step is to substitute the new found factors into the expression in the braces $\{ \}$ of equation (A V-1) and write

$$(x+a)(x+b)x+c=0$$

$$\frac{1}{c}(x+a)(x+b)x=-1$$

See figure A V-2.

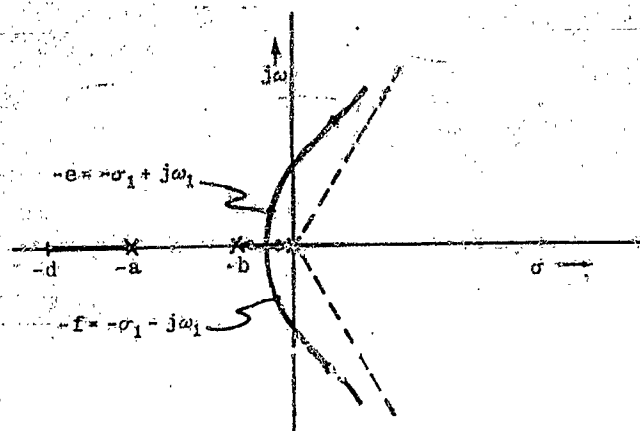


Figure A V-2. $[(x+A)x+B]x+C=0=(x+d)(x+e)(x+f)$

Finally, the factors $(x+d)$, $(x+e)$, and $(x+f)$ are substituted into equation (A V-1) and the new equation solved:

$$(x+d)(x+e)(x+f)x+D=0$$

$$\frac{1}{D}(x+d)(x+e)(x+f)x=-1$$

See Figure A V-3.

Evidently, this procedure can be used to solve for the roots of any rational polynomial with constant coefficients. The essence of the method is to reduce the polynomial to a series of equations of the form $1+P(x)=0$, any of the methods discussed in Chapter III may be used to handle it. In particular the open-loop/closed-loop method yields a first approximation to the roots that can be of great value. (See Section III-3.)

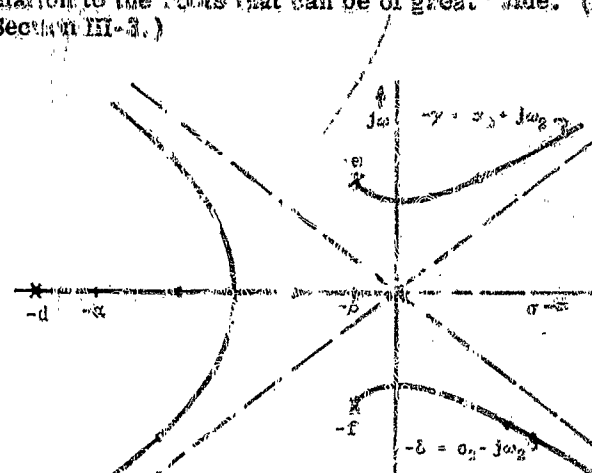


Figure A V-3. $[(x+A)x+B]x+Cx+D=0=x^3+A_1x^2+Bx+Cx+D$
 $=0=(x+d)(x+e)(x+f)x+D$